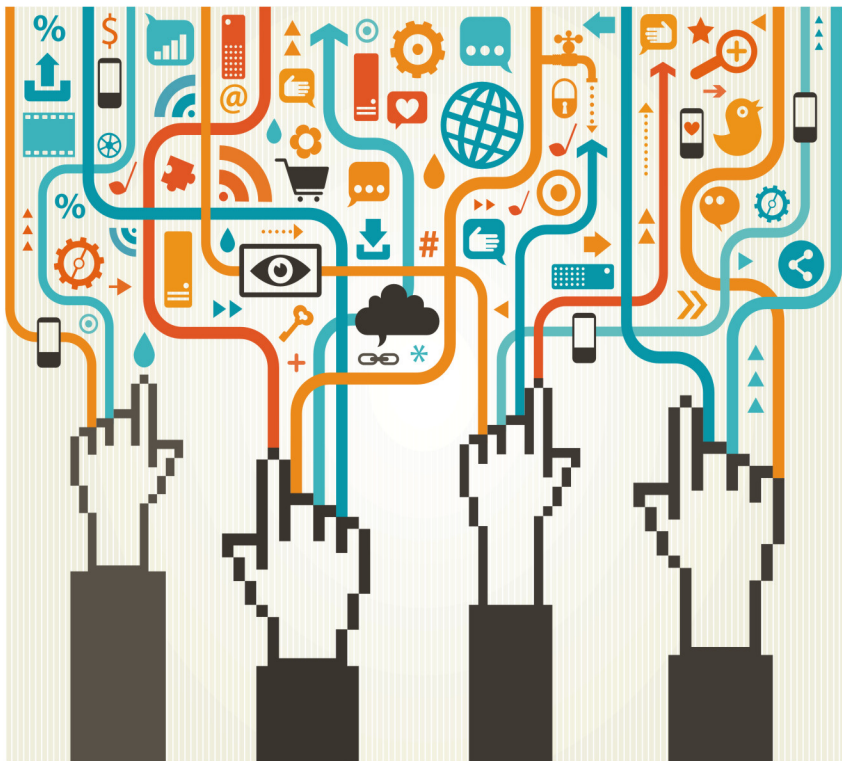


**Sociologie**  
**Antropologie**

**Marian Vasile**

# INTRODUCERE ÎN SPSS

**pentru cercetarea socială și de piață**



Collegium

# POLIROM

*Collegium. Sociologie. Antropologie* – serie coordonată de Lazăr Vlăsceanu  
și Liviu Chelcea și inițiată de Elisabeta Stănciulescu.

© 2014 by Editura POLIROM

Această carte este protejată prin copyright. Reproducerea integrală sau parțială, multiplicarea prin orice mijloace și sub orice formă, cum ar fi xeroxarea, scanarea, transpunerea în format electronic sau audio, punerea la dispoziția publică, inclusiv prin internet sau prin rețele de calculatoare, stocarea permanentă sau temporară pe dispozitive sau sisteme cu posibilitatea recuperării informațiilor, cu scop comercial sau gratuit, precum și alte fapte similare săvârșite fără permisiunea scrisă a deținătorului copyrightului reprezintă o încălcare a legislației cu privire la protecția proprietății intelectuale și se pedepsesc penal și/sau civil în conformitate cu legile în vigoare.

Reprints on pages: 33, 35, 38, 42-44, 46-50, 53, 56, 59-60, 64-65, 68, 72-73, 75-76, 81, 85, 87-88, 92, 99-100, 104-105, 108, 116-117, 119-120, 126-127, 129, 131-132, 139-140, 143-144, 146, 148, 150, 157, 160, 162-165, 167-170, 176, 179, 189-191, 196:  
Courtesy of International Business Machines Corporation, © International Business Machines Corporation

IBM, the IBM logo, ibm.com, and SPSS are trademarks or registered trademarks of **International Business Machines Corporation**, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "IBM Copyright and trademark information" at [www.ibm.com/legal/copytrade.shtml](http://www.ibm.com/legal/copytrade.shtml).

Pe copertă: © iStockphoto.com/DrAfter123

[www.polirom.ro](http://www.polirom.ro)

Editura POLIROM

Iași, B-dul Carol I nr. 4; P.O. BOX 266, 700506  
București, Splaiul Unirii nr. 6, bl. B3A, sc. 1, et. 1,  
sector 4, 040031, O.P. 53, C.P. 15-728

Descrierea CIP a Bibliotecii Naționale a României:

VASILE, MARIAN

*Introducere în SPSS pentru cercetarea socială și de piață: o perspectivă aplicată /*

Marian Vasile. – Iași: Polirom, 2014

Bibliogr.

ISBN print: 978-973-46-4136-9

ISBN ePub: 978-973-46-4460-5

ISBN PDF: 978-973-46-4462-9

004.42 SPSS:339.13+316.65

519.22:339.13+316.65

Printed in ROMANIA

**Marian Vasile**

# **INTRODUCERE ÎN SPSS**

**pentru cercetarea socială și de piață**

**O perspectivă aplicată**

POLIROM  
2014

MARIAN VASILE este lector universitar la Universitatea din București, Facultatea de Sociologie și Asistență Socială, și cercetător la Academia Română, Institutul de Cercetare a Calității Vieții. A încheiat un proiect postdoctoral centrat pe analiza mecanismelor care influențează fericirea indivizilor. Este doctor în sociologie din 2010, cu o teză în care identifică o serie de stiluri de viață specifice României anilor 2000. A participat la numeroase cursuri despre regresia multinivel și modele de ecuații structurale și la proiecte de cercetare socială în domenii precum stilurile de viață, stratificarea socială, educația, calitatea vieții și valorile. Este membru al European Sociological Association și International Sociological Association. Mai multe informații pe <http://www.marian-vasile.ro>.

# Cuprins

|   |          |
|---|----------|
| <i>Lista figurilor și tabelelor .....</i>   | <i>7</i> |
| 1. Introducere.....   | 11       |
| 1.1. Cui este utilă această carte.....  | 11       |
| 1.2. Ce conține această carte și cum să o citim.....                                      | 13       |
| 1.3. Materiale suplimentare .....   | 14       |
| 1.4. Mulțumiri .....  | 14       |
| 2. Crearea unei baze de date .....  | 17       |
| 2.1. În ce program introducem chestionarele ? .....                                       | 19       |
| 2.2. Cum ducem datele în SPSS ? .....   | 32       |
| 2.3. Exerciții .....  | 39       |
| 3. Gestionarea bazei de date .....  | 41       |
| 3.1. Câteva setări elementare (Edit > Options) .....                                      | 41       |
| 3.2. Pe scurt, despre structura programului :<br>Data și Variable View.....               | 49       |
| 3.3. Ponderarea bazei de date (Data > Weight Cases) .....                                 | 54       |
| 3.4. Filtrarea bazei de date (Data > Select Cases) .....                                  | 57       |
| 3.5. Separarea bazei de date (Data > Split File) .....                                    | 62       |
| 3.6. Exerciții .....  | 69       |
| 4. Curățarea și validarea unei baze de date.....  | 71       |
| 4.1. Etichetarea variabilelor și a valorilor variabilelor.....                            | 72       |
| 4.2. Definirea nonrăspunsurilor .....   | 79       |
| 4.3. Verificarea introducerii eronate a unor coduri.....                                  | 83       |
| 4.4. Validarea logică prin urmărirea filtrelor<br>și a unor întrebări factuale .....      | 84       |
| 4.5. Construirea unor variabile noi.....  | 92       |
| 4.6. Exerciții .....  | 94       |
| 5. Gestionarea variabilelor .....   | 95       |
| 5.1. Crearea unei alte variabile utilizând meniul<br>Recode into Different Variable ..... | 95       |
| 5.2. Crearea unei alte variabile utilizând meniul Compute.....                            | 103      |
| 5.3. Exerciții .....  | 108      |
| 6. O primă privire asupra datelor .....   | 111      |
| 6.1. Cum gândește majoritatea și cât de omogene<br>sunt grupurile comparate.....          | 112      |

|  |     |
|--|-----|
| 6.2. Asocierea dintre variabile categoriale.<br>Tabelul de contingență (Crosstabs) ..... | 123 |
| 6.3. Diferențe între medii : testul t pentru eșantioane<br>independente și ANOVA.....    | 138 |
| 6.4. Două grafice uzuale în descrierea datelor .....                                     | 145 |
| 6.5. Exerciții .....   | 151 |
| 7. Explorarea datelor : asumptii.....  | 153 |
| 7.1. Distribuția unei variabile .....  | 154 |
| 7.2. Relația liniară dintre două variabile .....   | 166 |
| 7.3. Soluții la încălcarea asumptiei de normalitate a distribuției .....                 | 171 |
| 7.4. Exerciții .....   | 172 |
| 8. Corelația și regresia liniară multiplă .....  | 173 |
| 8.1. Corelația liniară .....   | 174 |
| 8.2. Regresia liniară multiplă .....   | 185 |
| 8.3. Exerciții .....   | 201 |
| <i>Bibliografie</i> .....  | 203 |

## Lista figurilor și tabelelor

|  |    |
|--|----|
| Figura 2.1. Crearea unei baze de date în Excel. Pasul 1 :<br>introducerea denumirilor variabilelor .....                               | 21 |
| Figura 2.2. Impunerea condițiilor de introducere a datelor în Excel.<br>Tabul Settings : fereastra inițială .....                      | 21 |
| Figura 2.3. Impunerea condițiilor de introducere a datelor în Excel.<br>Tabul Settings : fereastra cu condiții .....                   | 22 |
| Figura 2.4. Impunerea condițiilor de introducere a datelor în Excel.<br>Tabul Error Alert : fereastra cu mesaj în stilul Warning ..... | 23 |
| Figura 2.5. Formular de introducere a datelor în Excel. Formular gol<br>și formular cu informații introduse pentru un respondent ..... | 24 |
| Figura 2.6. Bază de date creată în Excel .....   | 24 |
| Figura 2.7. Importarea datelor din formatul Excel în formatul SPSS.<br>Selectarea fișierului Excel .....                               | 33 |
| Figura 2.8. Importarea datelor din formatul Excel în formatul SPSS.<br>Selectarea foii de lucru care conține datele .....              | 33 |
| Figura 2.9. Ordonarea cazurilor : după o variabilă, de la valorile mici<br>la valorile mari ale acesteia .....                         | 35 |
| Figura 2.10. Unirea a două baze cu aceiași respondenți și variabile diferite .....   | 35 |
| Figura 2.11. Unirea a două baze cu respondenți diferiți și aceleași variabile .....  | 37 |
| Figura 2.12. Variabile cu proprietăți diferite (coloana Width din Variable View) .....   | 38 |
| Figura 3.1. Setări care cresc ușurința de utilizare a programului. Tabul General .....   | 42 |
| Figura 3.2. Tabul Output Labels (Edit > Options) : două tipuri de vizualizare<br>în Output .....                                       | 44 |
| Figura 3.3. Tabul Pivot Tables (Edit > Options) : modificarea<br>designului tabelor .....  | 46 |
| Figura 3.4. Tabul File Locations (Edit > Options) : fișierul de lucru<br>și jurnalul SPSS .....  | 47 |
| Figura 3.5. Tabul Syntax Editor (Edit > Options) : cum putem<br>face sintaxa mai ușor de utilizat .....                                | 48 |
| Figura 3.6. Data View .....  | 50 |
| Figura 3.7. Meniul Window > Split : rezultatul împărțirii .....  | 53 |
| Figura 3.8. Variable View .....  | 53 |
| Figura 3.9. Find : căutare după un cuvânt-cheie .....  | 56 |
| Figura 3.10. Data > Weight Cases : meniul în care activăm ponderea .....   | 56 |
| Figura 3.11. Confirmare vizuală că ponderea este activă .....  | 57 |
| Figura 3.12. Meniul Data > Select Cases : fereastra inițială prin care activăm,<br>dezactivăm, copiem sau ștergem cazuri .....         | 59 |

|  |     |
|--|-----|
| Figura 3.13. Meniul Data > Select Cases : fereastra cu filtrul<br>care menține active doar anumite cazuri .....    | 64  |
| Figura 3.14. Meniul Data > Select Cases : fereastra cu filtrul<br>care creează o bază nouă de date .....           | 65  |
| Figura 3.15. Meniul Data > Split File .....  | 68  |
| Figura 3.16. Meniul Analyze > Descriptive Statistics > Frequencies :<br>cum calculăm media unei variabile .....    | 68  |
| Figura 4.1. Variable View : baza de date înainte și după etichetare .....  | 72  |
| Figura 4.2. Fișierul de sintaxă : afișarea listei derulante de comenzi .....                                       | 73  |
| Figura 4.3. Meniul Utilities > Variables : cum găsim rapid o variabilă<br>și cum îi copiem numele în sintaxă ..... | 75  |
| Figura 4.4. Variable View : etichetarea valorilor variabilei .....   | 76  |
| Figura 4.5. Definirea nonrăspunsurilor în Variable View .....  | 81  |
| Figura 4.6. Crosstabs : realizarea unui tabel de contingență .....   | 85  |
| Figura 4.7. Recode into Same Variables .....   | 87  |
| Figura 4.8. Crearea de noi variabile (Compute) .....   | 92  |
| Figura 5.1. Meniul Transform > Recode into Different Variables .....   | 99  |
| Figura 5.2. Transform > Compute : crearea unei variabile de ponderare .....  | 104 |
| Figura 5.3. Relație nonliniară dintre vârstă și satisfacția cu viața .....   | 106 |
| Figura 5.4. Distribuția venitului înainte și după logaritmare .....  | 107 |
| Figura 5.5. Realizarea graficelor din meniul Frequencies > Charts .....  | 108 |
| Figura 6.1. Meniul Frequencies, butonul Statistics .....   | 116 |
| Figura 6.2. Meniul Descriptives .....  | 119 |
| Figura 6.3. Meniul Explore .....   | 120 |
| Figura 6.4. Meniul Crosstabs .....   | 126 |
| Figura 6.5. Grafic bară obținut folosind meniul Crosstabs .....  | 128 |
| Figura 6.6. Editarea unui tabel de contingență în Output (Pivot) .....   | 131 |
| Figura 6.7. Meniul Independent-Samples T Test .....  | 139 |
| Figura 6.8. Meniul One-Way ANOVA. Analiza de varianță .....  | 143 |
| Figura 6.9. Meniul Frequencies, Charts .....   | 146 |
| Figura 6.10. Grafic bară .....   | 147 |
| Figura 6.11. Editarea valorilor de pe bare, Properties .....   | 148 |
| Figura 6.12. Histograme .....  | 149 |
| Figura 6.13. Chart Editor, Properties pentru histogramă .....  | 150 |
| Figura 7.1. Grafic bară pentru verificarea asumției de normalitate<br>a distribuției .....                         | 156 |
| Figura 7.2. Meniul Explore .....   | 157 |
| Figura 7.3. Box-plot : distribuția satisfacției cu viața în funcție<br>de poziția socială subiectivă .....         | 158 |
| Figura 7.4. Find .....   | 160 |
| Figura 7.5. Meniul Window > Split .....  | 160 |
| Figura 7.6. Normal probability plot .....  | 161 |
| Figura 7.7. Normal P-P Plots .....   | 162 |
| Figura 7.8. Descriptives, scoruri z .....  | 162 |
| Figura 7.9. Find variables .....   | 163 |



|  |     |
|--|-----|
| Figura 7.10. Sort Cases.....   | 163 |
| Figura 7.11. Meniurile Frequencies și Descriptives : calcularea skewness<br>și kurtosis .....                                | 164 |
| Figura 7.12. Scatterplot (nor de puncte) care arată o relație perfect liniară.....   | 167 |
| Figura 7.13. Meniul Graphs > Chart Builder .....   | 168 |
| Figura 7.14. Window > Split : consultarea vizuală în Data View<br>a unor inadvertențe în date.....                           | 170 |
| Figura 7.15. Limitele scatterplotului : când o variabilă are puține valori.....  | 171 |
| Figura 8.1. Relația nonliniară dintre vârstă și satisfacția cu viața.....  | 174 |
| Figura 8.2. Meniul corelației bivariate (Correlate > Bivariate).....   | 176 |
| Figura 8.3. Meniul corelației parțiale (Correlate > Partial).....  | 179 |
| Figura 8.4. Meniul Analyze > Regression > Linear .....   | 189 |
| Figura 8.5. Scatterplot care ne arată un caz extrem .....  | 195 |
| Figura 8.6. Meniul Simple Scatterplot.....   | 196 |
|  |     |
| Tabelul 2.1. Două tipuri de întrebări folosite în chestionare : închise și deschise ...                                      | 17  |
| Tabelul 3.1. Meniuri frecvent utilizate.....   | 51  |
| Tabelul 3.2. Tabel de frecvență : Care sunt codurile folosite pentru bărbați<br>și pentru femei ? .....                      | 58  |
| Tabelul 3.3. Tabel de frecvență : verificarea corectitudinii filtrului.....  | 60  |
| Tabelul 3.4. Tabel de frecvență : Distribuția fericirii<br>în rândul bărbaților români (WVS 2012) .....                      | 61  |
| Tabelul 3.5. Tabel de frecvență : Care sunt codurile pentru bărbați<br>și pentru mediul rural ?.....                         | 62  |
| Tabelul 3.6. Tabele de frecvență : verificarea corectitudinii<br>filtrului V240 = 1 & (V248 = 8   V248 = 9) .....            | 63  |
| Tabelul 3.7. Tabel de frecvență : Distribuția fericirii<br>în rândul bărbaților români cu studii superioare (WVS 2012) ..... | 63  |
| Tabelul 3.8. Reprezentare grafică a rezultatului separării bazei de date .....   | 67  |
| Tabelul 3.9. Media vârstei : tabel obținut prin separarea bazei de date .....  | 69  |
| Tabelul 4.1. Tabel de frecvență înainte de definirea nonrăspunsurilor.....   | 82  |
| Tabelul 4.2. Tabel de contingență ce verifică un filtru, dar este folosit<br>pentru validare logică (1) .....                | 84  |
| Tabelul 4.3. Tabel de contingență care verifică un filtru, dar este folosit<br>și pentru validare logică (2) .....           | 89  |
| Tabelul 4.4. Validare logică : tabel de contingență .....  | 91  |
| Tabelul 5.1. Niveluri de măsurare .....  | 97  |
| Tabelul 5.2. Tabel de frecvență pentru autoevaluarea stării de sănătate .....  | 101 |
| Tabelul 5.3. Tabel de contingență pentru verificarea corectitudinii recodificării ....                                       | 102 |
| Tabelul 6.1. Tabel de frecvență : înainte și după definirea nonrăspunsurilor.....  | 114 |
| Tabelul 6.2. Tabel de frecvență : după definirea nonrăspunsurilor .....  | 115 |
| Tabelul 6.3. Tabele de frecvență și indicatori statistici ai tendinței<br>centrale și ai variației .....                     | 118 |

|  |     |
|--|-----|
| Tabelul 6.4. Tabel Outliers obținut din meniul Explore .....   | 121 |
| Tabelul 6.5. Output produs de meniul Explore .....   | 122 |
| Tabelul 6.6. Tabele de frecvență : inspectarea variabilelor înainte<br>de analiza de contingență (Crosstabs) .....                         | 125 |
| Tabelul 6.7. Tabel de contingență care conține doar frecvențe<br>absolute (Count) .....  | 129 |
| Tabelul 6.8. Tabel de contingență care conține frecvențe absolute<br>și procente pe rând .....   | 130 |
| Tabelul 6.9. Testul Pearson chi-square : valoare și p .....  | 135 |
| Tabelul 6.10. Tabel de contingență cu statistici : stare civilă și fericire .....  | 136 |
| Tabelul 6.11. Tabel de contingență cu statistici : stare civilă<br>și încredere în oameni .....  | 137 |
| Tabelul 6.12. Testul t pentru eșantioane independente : output .....   | 141 |
| Tabelul 6.13. Rezultate ale analizei de varianță .....   | 144 |
| Tabelul 7.1. Tabel de frecvență : verificarea variației variabilelor categoriale ....  | 154 |
| Tabelul 7.2. Skewness și kurtosis. Calcule efectuate în meniul Explore.<br>Tabele obținute prin pivotare .....                             | 165 |
| Tabelul 8.1. Output cu sau fără opțiunile Flag... sau Display... în meniurile<br>corelației bivariate, respectiv corelației parțiale ..... | 181 |
| Tabelul 8.2. Opțiunea Means and standard deviations din meniurile corelației<br>bivariate, respectiv corelației parțiale .....             | 182 |
| Tabelul 8.3. Opțiunea Zero-order correlations în meniul corelației parțiale .....  | 183 |
| Tabelul 8.4. Output regresie liniară multiplă, Descriptives statistics .....   | 198 |
| Tabelul 8.5. Output regresie liniară, Corelații bivariate .....  | 198 |
| Tabelul 8.6. Output regresie liniară, R2 .....   | 199 |
| Tabelul 8.7. Output regresie liniară, Coefficients .....   | 200 |

# 1. Introducere

## 1.1. Cui este utilă această carte

De câțiva ani buni, îi ajut pe studenți să înțeleagă utilitatea statisticii în cercetarea socială și de marketing. Implicit, pentru că realizarea analizelor statistice fără un software dedicat este dificil de imaginat astăzi, încerc să îi familiarizez cu unul dintre acestea. Pentru că studenții cu care lucrez sunt, într-un număr destul de mare, absolvenți de filologie sau de științe sociale, acest demers este o provocare, una plăcută însă. Programul de statistică utilizat în această lucrare este **IBM® SPSS® Statistics software (SPSS)**<sup>1</sup>, versiunea 17. Toate operațiunile pot fi reproduse folosind orice versiune recentă a programului. Vă recomand versiunile mai noi, pentru că sintaxa reliefează, folosind culori, diferitele elemente care o compun. Va fi mai ușor să vă obișnuiți cu aceasta.

Studenții optează pentru un curs doar dacă acesta li se pare util. Statistica este cât se poate de utilă în orice domeniu, dar în științele sociale este destul de greu să îi convingi pe cei care se tem de matematică să aleagă de bunăvoie și nesiliți de nimeni să treacă prin acest „calvar”. Unii studenți procedează în felul următor : deschid programul de statistică, în cazul de față SPSS, și încearcă să reproducă pașii explicați în diferite manuale sau tutoriale. Inevitabil, interacționează cu concepte din statistică, dar le acordă mai puțină importanță în procesul de învățare decât meniurilor și comenzilor din program. Aceasta este o perspectivă „inversă”. Nu poți învăța să folosești un program de statistică dacă nu știi... statistică. Este ca și când ai vrea să devii pilot de Formula 1 fără să ai permis de conducere. Această abordare duce la învățare mecanică : utilizatorul intră în meniurile SPSS și dă clickuri ici și colo fără să-i fie clar de ce face aceste lucruri, de ce alege o opțiune, și nu alta, sau cum sunt interpretate rezultatele pe care le oferă aceste acțiuni.

Cum ar trebui să procedeze studentul ? Ar trebui să parcurgă un manual de statistică și, simultan, un manual în care analizele statistice sunt puse în practică într-un program de statistică – SPSS. Cursurile de statistică teoretică sau aplicată nu sunt niciodată suficiente. Domeniul este atât de dezvoltat, încât subiectul nu poate fi epuizat într-o singură lucrare. Învățarea statisticii este un proces. Cel sau cea care se angajează în acest demers trebuie să adauge consultării materialelor teoretice multe exerciții folosind date reale. Astăzi este foarte la îndemână acest lucru. Tot mai

---

1. SPSS Inc. a fost achiziționat de IBM în octombrie 2009.

multe date sunt accesibile gratuit. Vedeți în acest sens studiile European Values Study, World Values Survey, European Social Survey, Eurobarometrul etc. Pe paginile web ale acestor cercetări găsiți chestionarele utilizate, documentație extensivă despre activitatea de teren, baze de date și multe alte informații care vă ajută să înțelegeți complexitatea abordării cantitative a realității și tipul de rezultate care pot fi obținute astfel. Să presupunem că Maria și-a dat seama că statistica este importantă și s-a decis să învețe principalele tehnici utilizate în piață. Dar este abia în anul I de facultate, astfel că nu a avut ocazia să participe la cercetări în calitate de analist. Adică nu i-a pus nimeni în brațe un chestionar, o bază de date și o listă de întrebări de cercetare pentru soluționarea cărora să fie nevoită să facă anumite analize statistice. În această situație, ar putea să rezolve exercițiile din manualele de statistică folosind, evident, programul SPSS. Din experiența proprie, pot spune că, într-un final, va ajunge să înțeleagă multe lucruri, dar pe parcurs s-ar putea să se descurajeze și să aibă impresia că drumul pe care s-a angajat este foarte greu și nu tocmai plăcut. Dacă nu este autodidactă sau foarte hotărâtă, atunci Maria s-ar putea să renunțe la un moment dat.

Consider că lipsește un manual care să îl ajute pe studentul începător în cercetare să unească logica activității de cercetare în științele sociale și logica manualelor de statistică. Principalele întrebări la care răspunde acest volum sunt :

- Care este legătura dintre chestionarul care a fost utilizat pentru a culege date și baza de date ?
- Cum realizați o bază de date ?
- Ce înseamnă să curățați baza de date ?
- Ce înseamnă să pregătiți datele pentru analiză ?
- Ce sunt codificarea și recodificarea unei variabile ?
- Cum creați variabile într-o bază de date ?
- De ce trebuie să vă uitați la date, înainte de a face analiza care vă interesează ?
- Cum faceți această explorare primară a datelor ?
- Ce este un tabel ? Dar un tabel de contingență ?
- Cum verificați dacă variabila X este asociată cu variabila Y ?
- Care este diferența dintre asociere și corelație ?
- Dacă doriți să explicați un fenomen, să zicem fericirea (Y), iar teoriile vă spun că este posibil ca acesta să fie explicat de mai mulți factori, să zicem starea de sănătate (X1), calitatea relațiilor sociale (X2) și cantitatea de timp liber avută la dispoziție (X3), ce tehnică statistică puteți folosi în acest sens ?

Lista nu este completă. Pe măsură ce citiți acest volum, puteți adăuga întrebările la care ați găsit un răspuns. Statistica oferă mai multe metode prin care putem răspunde la aceeași întrebare. SPSS oferă mai multe comenzi pentru aceeași analiză. Le voi prezenta pe cele mai importante pentru cei aflați la început de drum. Tranziția spre lucrurile mai dificile va fi mai ușoară după ce ați parcurs acest volum.

*Introducere în SPSS pentru cercetarea socială și de piață. O perspectivă aplicată* se adresează, în primul rând, studenților care vor să facă primii pași în abordarea cantitativă a socialului. Ei pot fi studenți la sociologie, marketing sau

administrarea afacerilor. Logica este în multe situații similară în aceste domenii. Apoi, sunt vizați masteranzii care au o pregătire limitată în statistică și utilizarea SPSS-ului, dar și doctoranzii care nu au urmat un curs intensiv în acest domeniu și nici nu au lucrat în multe proiecte care folosesc date cantitative. De asemenea, cred că este util și pentru cercetătorii care au utilizat SPSS, dar l-au învățat, mai degrabă, „pe încercate”, și nu în mod sistematic.

## 1.2. Ce conține această carte și cum să o citim

Acord o atenție considerabilă aspectelor premergătoare activității de analiză cantitativă a datelor culese prin aplicarea unor chestionare. O mare parte din timpul activității de analiză este consumat de aceste aspecte preliminare. Primele elemente care îl preocupă pe cercetătorul cantitativist sunt elaborarea bazei de date (capitolul 2) și curățarea acesteia (capitolul 4). Pentru procesul de curățare, acesta trebuie să învețe câteva operațiuni cum ar fi filtrarea bazei de date (capitolul 3) sau crearea de variabile noi (capitolul 5). Este dificil să scrii o lucrare care urmărește toți acești pași, exact în ordinea în care se întâmplă în realitate. Demersul este circular, de aceea, de exemplu, în procesul de curățare voi folosi informații prezentate și în capitolele ulterioare, cum ar fi cele despre tabelele de contingență (capitolul 6). Cert este că informațiile din capitolele 2, 3, 4 și 5 sunt esențiale și trebuie citite înainte de a trece la capitolul 6. Odată cu capitolul 6, cititorul primește și informații despre analizele statistice uzuale care pot fi utilizate pentru a răspunde la întrebări de cercetare. Cum observăm modul în care gândește majoritatea? Cât de omogene sunt diferite grupuri în funcție de o anumită caracteristică? Media, mediana, abaterea standard și altele sunt doar câteva elemente utile pentru a răspunde la astfel de întrebări. Tabelul de contingență ne va ajuta să vedem dacă două variabile sunt independente sau nu. Apoi, aflăm cum putem testa diferența dintre două sau mai multe grupuri în funcție de o caracteristică. După aceea, aflăm cum explicăm variația unei variabile în funcție de mai multe caracteristici. Media generală la învățatură a elevilor care au făcut trei ore de educație fizică pe săptămână la școală este mai ridicată decât cea a elevilor care au făcut cel mult o oră de educație fizică pe săptămână la școală? Volumul vânzărilor iaurtului cu căpsuni produs de firma „Iaurt pentru toți” este mai mare dacă în hipermarketuri se folosește testarea produsului (adică firma a angajat promotori care le oferă potențialilor cumpărători să guste iaurtul respectiv) decât dacă nu se folosește? Informațiile prezentate în acest volum pot fi utilizate atât în situații întâlnite în cercetarea socială, cât și în cea de piață. La acest gen de întrebări putem răspunde statistic folosind informațiile din capitolul 6. Realitatea socială este mult mai complexă. Nu ne putem aștepta ca media generală la învățatură a elevilor să depindă doar de practicarea frecventă a unor activități sportive, la fel cum nu ne putem aștepta ca volumul vânzărilor unui tip de iaurt

să depindă doar de prezența promotorilor în magazine. Am putea adăuga, pentru primul caz, numărul de ore petrecute în bibliotecă studiind individual, ajutorul primit din partea părinților, participarea la activități extrașcolare cu caracter educativ, numărul colegilor sau prietenilor cu care elevul își petrece timpul liber, caracteristicile acestora etc. În al doilea caz, am putea adăuga calitatea distribuției, atractivitatea ambalajului, prețul produsului, poziționarea la raft etc. Avem, așadar, o variabilă dependentă și mai multe variabile independente. Pentru acest gen de situații, informațiile prezentate în capitolul 8 vor fi utile.

Capitolul 7 tratează o serie de asumptii fundamentale pentru analizele statistice prezentate în capitolele 6 și 8. Aș fi putut opta pentru o prezentare sumară în cadrul fiecărui capitol, însă am vrut să subliniez importanța acestui pas. Domeniul explorării asumptiilor este vast, depășind obiectivele acestui volum care constituie, în primul rând, un material introductiv. Pe măsură ce învățați mai multe analize statistice, în special multivariate, veți identifica și alte asumptii care trebuie testate.

Toate capitolele se încheie cu o listă de exerciții care pot fi folosite pentru a pune în practică informațiile prezentate pe parcursul capitolului respectiv. Exercițiile înseamnă experiență acumulată atât cu conceptele, cât și cu programul de statistică. Consider că niciodată nu facem suficiente exerciții, așadar lista cu exerciții de la finalul fiecărui capitol este doar un prolog al eforturilor dumneavoastră viitoare.

## 1.3. Materiale suplimentare

Puteți descărca date utilizate pentru diferite exemple, sintaxe și outputuri produse de aceste sintaxe de pe pagina de internet: <http://www.marian-vasile.ro/publications/spss>.

## 1.4. Mulțumiri

Aș vrea să le mulțumesc celor de la care am învățat, la rândul meu, multe dintre lucrurile pe care le știu atât despre analizele statistice, cât și despre utilizarea SPSS. În primul rând, vreau să îi mulțumesc lui Bogdan Voicu, care nu numai că mi-a răspuns la toate întrebările, dar mi-a oferit și oportunitatea de a-i fi asistent la cursurile sau trainingurile ținute în diferite contexte. Apoi, aș vrea să îi mulțumesc lui Alexandru Cernat pentru că a acordat timp citirii acestui material, oferindu-mi sugestii extrem de utile. Îi mulțumesc lui Ioan Mărginean pentru că m-a provocat să gândesc critic diferite situații întâlnite în cercetarea calității vieții, și nu numai. Nu în ultimul rând, le mulțumesc lui Liviu Chelcea și Lazăr Vlăsceanu pentru că m-au încurajat să public această lucrare.

Multe dintre informațiile acumulate și transpuse, într-o formă sau alta, în acest volum au fost acumulate în cadrul unor proiecte de cercetare similare cu cel postdoctoral susținut de UEFISCDI, care s-a derulat între 2011 și 2013 sub denumirea *Drumuri diferite către o viață mai bună: comparații internaționale longitudinale ale determinanților satisfacției cu viața* (PN-II-RU-PD-2011-3-0117). Un alt proiect este cel coordonat de Bogdan Voicu, care s-a derulat între 2011 și 2014 sub titlul *Schimbarea socială în contextul migrației internaționale: patternuri valorice, participare civică și politică, satisfacția cu viața* (PN-II-ID-PCE-2011-3-0210). Pentru mai multe detalii, puteți consulta paginile de internet <http://www.stilurideviata.ro> și <http://www.romanianvalues.ro>.





## 2. Crearea unei baze de date

„Cercetare cantitativă” sau „analiză cantitativă” sunt două concepte frecvent folosite de practicieni în activitatea de zi cu zi. Ambele fac trimitere la culegerea și analiza unor informații prin utilizarea chestionarului ca instrument de cercetare. Chestionarul cuprinde o serie de întrebări închise și, uneori, și întrebări deschise. Întrebările închise au răspunsurile predefinite de cercetător, persoana care este rugată să răspundă la întrebare (respondentul) trebuind doar să îl aleagă pe cel care i se potrivește cel mai bine. Întrebările deschise nu au răspunsuri predefinite, respondentul trebuind să compună, folosind cuvintele proprii, un răspuns care caracterizează cel mai bine modul cum gândește, se comportă sau, mai general, care prezintă situația sa la momentul interviuării sau la cel de referință folosit de cercetător. Tabelul 2.1 prezintă un exemplu care diferențiază aceste două tipuri de întrebări.

**Tabelul 2.1.** Două tipuri de întrebări folosite în chestionare : închise și deschise

| Întrebare închisă  | Întrebare deschisă   |
|--|--|
| D2. Ocupația dvs. actuală (principală) :<br>1. agricultor<br>2. muncitor (meseriaș)<br>3. tehnician, maistru, funcționar<br>4. ocupații cu studii superioare<br>5. altă ocupație<br>6. elev, student<br>7. pensionar<br>8. casnică<br>9. acum sunt șomer<br>10. patron | Q112. Cum se numește munca pe care o desfășurați (la principalul loc de muncă) ?<br><br>Q112a. Ce fel de activitate desfășurați în cea mai mare parte a timpului ? |
| Sursa : <i>Diagnoza calității vieții din România</i> , Institutul de Cercetare a Calității Vieții, 2010.   | Sursa : <i>European Values Study</i> , Institutul de Cercetare a Calității Vieții, 2008.   |

În exemplul din tabelul 2.1, cercetătorul este interesat să afle structura ocupațiilor din România. Dacă folosește întrebarea închisă, atunci respondentul va alege varianta de răspuns care se potrivește cel mai bine situației sale. Dacă folosește întrebarea deschisă, atunci respondentul va descrie în cuvinte, cât mai detaliat, situația sa cu privire la acest subiect. De regulă, preferăm să folosim întrebări închise în chestionare pentru că aplicarea acestora durează mai puțin, sunt mai

ușor de înțeles, se introduc mai repede în baza de date, ne reprezentăm mai ușor ce fel de analize statistice putem realiza cu ele etc. Dacă optează pentru această variantă, cercetătorul trebuie să se asigure că lista variantelor de răspuns este completă, iar acestea nu se suprapun, adică respondentul nu se poate regăsi în mai multe răspunsuri simultan. Există situații, însă, în care mai multe răspunsuri sunt plauzibile pentru aceeași persoană, acestea fiind surprinse prin întrebările cu răspuns multiplu. Dacă optează pentru varianta cu întrebări deschise, atunci cercetătorul trebuie să știe cum va codifica răspunsurile primite. Codificarea presupune ca, din lista lungă de răspunsuri primite, să construiască una mai restrânsă, astfel încât fiecare categorie să poată primi un cod unic care va fi introdus în baza de date și, ulterior, va fi folosit pentru diferite analize statistice. De exemplu, întrebarea D2 (tabelul 2.1), are coduri de la 1 la 10. În acest exemplu, codificarea la Q112 și Q112a (tabelul 2.1), va fi realizată folosind o schemă de coduri standardizată, *International Standard Classification of Occupations*<sup>1</sup> (ISCO; Clasificarea internațională standard a ocupațiilor). Aceasta are mai multe variante. Dacă ne uităm la ISCO-88, putem vedea că, la nivelul cel mai înalt de generalitate, din răspunsurile deschise putem obține zece coduri, numerotate de la 0 la 9. Fiecare dintre aceste categorii ocupaționale largi poate fi divizată în mai multe subgrupuri. La cel mai rafinat nivel de specificare se poate ajunge la 390 de grupuri ocupaționale, adică 390 de coduri. Nivelul de detaliu ales de cercetător depinde, în mare măsură, de volumul eșantionului pe care îl folosește.

Informații despre tipurile de întrebări, regulile de elaborare a acestora, opțiunea pentru o formă sau alta și nu numai pot fi găsite în lucrările dedicate subiectului cum ar fi cele elaborate de Mărginean (1982), Bradburn, Sudman *et al.* (2004), Chelcea (2007), Saris și Gallhofer (2007) sau Malhotra (2007). Acestea sunt cunoștințe complementare celor prezentate aici și trebuie însușite pentru o înțelegere adecvată a procesului cercetării cantitative.

După aplicarea chestionarelor, acestea trebuie introduse în baza de date. Apoi baza de date trebuie curățată. Abia după aceste etape, putem trece la analizele statistice prin care răspundem la întrebările de cercetare. În acest capitol vom afla cum se realizează o bază de date în care sunt introduse chestionare și cum ajungem la baza de date în format SPSS. În capitolul 3 vom învăța câteva comenzi esențiale pentru gestionarea bazei de date, iar în capitolul 4 vom afla care sunt etapele procesului de curățare a bazei de date și ce presupune fiecare dintre ele.

Să presupunem că avem 1.000 de chestionare care trebuie introduse într-o bază de date. Pentru aceasta, există mai multe opțiuni. Astăzi, din ce în ce mai multe institute și companii de cercetare socială și/sau de piață înlocuiesc chestionarele pe hârtie cu chestionarele în format digital. Adică operatorul de teren nu

---

1. Documentele despre această clasificare pot fi consultate pe pagina dedicată de pe site-ul ISCO : <http://www.ilo.org/public/english/bureau/stat/isco/isco88/publ3.htm>.

mai completează cu pixul pe hârtie răspunsurile la întrebări, ci dă click pe un laptop sau pe o tabletă. Formatul digital de aplicare a chestionarelor are mai multe avantaje față de cel clasic, pe hârtie. Baza de date, chiar și cea în format SPSS, este creată direct, fiind redus astfel necesarul de resurse umane, timp și bani pentru finalizarea cercetării. De asemenea, numărul erorilor întâlnite în procesul de introducere a datelor este redus considerabil. Nu în ultimul rând, activitatea operatorului de teren poate fi mai bine controlată. Deși investiția inițială în tablete sau alte instrumente electronice care pot fi utilizate pentru aplicarea chestionarelor este costisitoare, pe termen lung, investiția se amortizează și își relevă utilitatea. Din ce în ce mai frecvent, chestionarele se aplică și on-line. Respondentul primește un link prin care poate accesa chestionarul pe care îl completează singur. Există o mulțime de soluții pentru această tehnică, cum ar fi și cea de tip *open source*, LimeSurvey<sup>1</sup>. Mai putem adăuga aplicarea prin telefon sau e-mail și, poate mai rar, prin poștă. Pentru detalii despre fiecare în parte, puteți consulta manuale de metodologie a cercetării sociale sau de piață cum ar fi *Marketing Research. An Applied Approach* (Malhotra și Birks, 2007).

Mulți studenți, masteranzi, doctoranzi, cercetători sau chiar firme și institute de cercetare nu își permit achiziționarea unor tablete cu software dedicat acestor acțiuni. De aceea, utilizează, în continuare, chestionarele tipărite pe hârtie care trebuie introduse într-o bază de date. Apoi această bază de date trebuie curățată. Soluții la îndemână în aceste situații sunt cele oferite de programele din suita Microsoft Office, mai exact, Microsoft Excel și Microsoft Access. În proiectele la care am lucrat, de cele mai multe ori, am introdus datele într-o bază de date realizată cu ajutorul Microsoft Access. În continuare, voi descrie pașii prin care realizăm o bază de date pentru introducerea chestionarelor folosindu-ne de Microsoft Excel, apoi de Microsoft Access. Apoi vom afla cum aducem în SPSS datele introduse într-unul dintre aceste programe.

## 2.1. În ce program introducem chestionarele ?

Acest subcapitol se referă la cercetările în care chestionarele sunt tipărite și aplicate de un operator de interviu prin procedeul față în față sau în care acestea sunt completate pe hârtie de către respondenți.

SPSS are propriile soluții de introducere a datelor. Mai multe detalii despre acestea și alte programe din domeniu pot fi găsite pe pagina de internet a producătorului programului<sup>2</sup>.

---

1. <https://www.limesurvey.org/en>.

2. <http://www-01.ibm.com/software/analytics/spss>.

### 2.1.1. Introducerea datelor în Microsoft Excel

Microsoft Excel este un program indispensabil în activitatea de cercetare, cu ajutorul căruia putem face diferite calcule, tabele sau grafice. Pe lângă acestea, poate fi folosit și pentru introducerea într-o bază de date a răspunsurilor la întrebările din chestionare. De fapt, vom alege un software sau un altul în funcție, în principal, de răspunsul la următoarea întrebare: pot fi evitate erorile de introducere a datelor? Altfel spus, dacă vrem să introducem răspunsurile la variabila gen, adică 1 = bărbat sau 2 = femeie, putem evita introducerea din greșeală a codului 3? Când introduceți unu-două chestionare, aceasta nu este o problemă deosebită pentru că puteți observa greșeala neintenționată. Dar dacă introduceți 300 de chestionare, fiecare având 100 de variabile, de la un moment dat nu mai observați la fel de ușor acest gen de greșeală. Dacă variabila dinaintea sau de după gen include printre codurile valide valoarea 3, atunci chestionarul poate fi introdus decalat. Vom vedea în capitolul 4, dedicat procesului de curățare, că astfel de erori pot fi identificate, dar corectarea lor presupune timp suplimentar de lucru. Așadar, dacă pot fi puse condiții care permit introducerea doar a codurilor corecte, atunci software-ul respectiv este adecvat pentru introducerea datelor. Un alt factor care contează în alegerea programului în care introducem datele constă în posibilitatea de a crea un formular de introducere care este plăcut privirii și care nu îl obosește pe operator.

Să presupunem că avem un chestionar cu trei variabile: id (o variabilă care are un cod unic pentru fiecare respondent), v1 (gen, unde 1 = bărbat sau 2 = femeie) și v2 (tipul de băutură carbogazoasă preferată, unde 1 = apă, 2 = suc cu cofeină sau 3 = suc de fructe). Vom prezenta în continuare o metodă rapidă de a elabora o bază de date în Microsoft Excel și de a introduce date în aceasta. De exemplu, la întrebări ne-au răspuns șase persoane. Așadar, trebuie să avem șase valori diferite la id. Pentru simplitate, acestea vor fi 1, 2, 3, 4, 5 și 6. Respondentul 1 este bărbat, deci primește codul 1. Acesta preferă apa, primind codul 1. Respondentul 6 este femeie, deci primește codul 2. Aceasta preferă sucul de fructe, primind codul 3. Datele care trebuie introduse sunt:

| id | v1 | v2 |
|----|----|----|
| 1  | 1  | 1  |
| 2  | 1  | 2  |
| 3  | 1  | 2  |
| 4  | 2  | 1  |
| 5  | 2  | 3  |
| 6  | 2  | 3  |

Deschidem o foaie goală în programul Microsoft Excel, iar cursorul ne duce în celula A1, adică la intersecția coloanei A cu rândul 1. Coloanele reprezintă variabilele (id, v1, v2), iar rândurile reprezintă respondenții (cele șase persoane). În celula A1 scriem id. În celula B1 scriem v1. În celula C1 scriem v2. Rezultatul ar trebui să arate ca în figura 2.1.

**Figura 2.1.** Crearea unei baze de date în Excel. Pasul 1 :  
introducerea denumirilor variabilelor

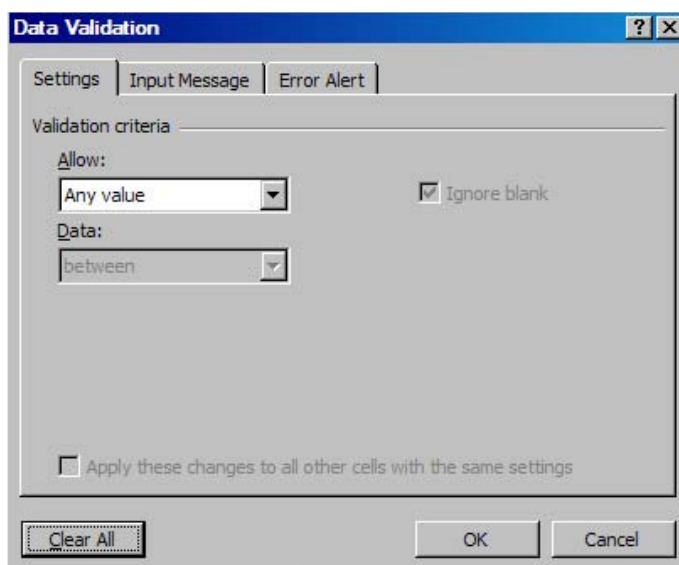
|   | A  | B  | C  |
|---|----|----|----|
| 1 | id | v1 | v2 |
| 2 |    |    |    |
| 3 |    |    |    |
| 4 |    |    |    |
| 5 |    |    |    |
| 6 |    |    |    |

Înainte de a introduce datele, trebuie să stabilim condițiile pentru fiecare variabilă :

- id să aibă valori cuprinse doar între 1 și 6 ;
- v1 să aibă doar valorile 1 sau 2 ;
- v2 să aibă valori cuprinse doar între 1 și 3.

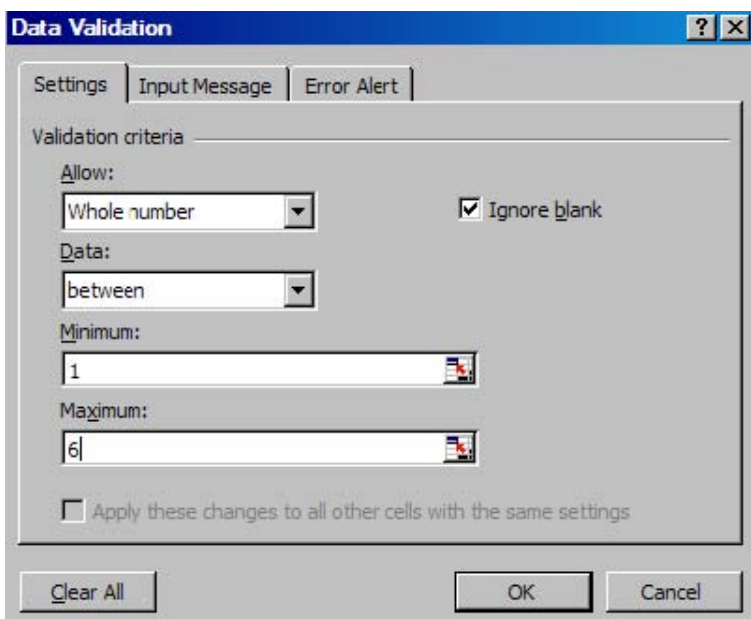
Ne ducem cu mouse-ul deasupra literei A, la numele primei coloane, și dăm click. Prin această operație, selectăm conținutul întregii coloane A. Apoi mergem în meniul **Data > Validation**. Se va deschide fereastra din figura 2.2. Ne interesează opțiunile din taburile **Settings** și **Error Alert**. În tabul **Settings** definim condiția. În tabul **Error Alert** vom scrie un mesaj de atenționare pentru operatorul care introduce chestionarele, precizând variantele corecte acceptate de celulele respective.

**Figura 2.2.** Impunerea condițiilor de introducere a datelor în Excel.  
Tabul Settings : fereastra inițială



Pentru că am selectat coloana A, adică variabila id, trebuie să impunem condițiile pentru aceasta : pot fi introduse doar valorile 1, 2, 3, 4, 5 sau 6. În tabul **Settings**, secțiunea **Allow** : selectăm **Whole number**. Se va activa secțiunea **Data** : în care lăsăm selectat **between**. Pentru că am lăsat selectat **between**, se activează alte două secțiuni, **Minimum** : și **Maximum** : , în care introducem codul 1, respectiv 6 (figura 2.3). Dacă operatorul introduce din greșală codul 7, nepermis în acest exemplu, atunci programul îl va avertiza că face o eroare înainte de a-i permite să continue introducerea datelor.

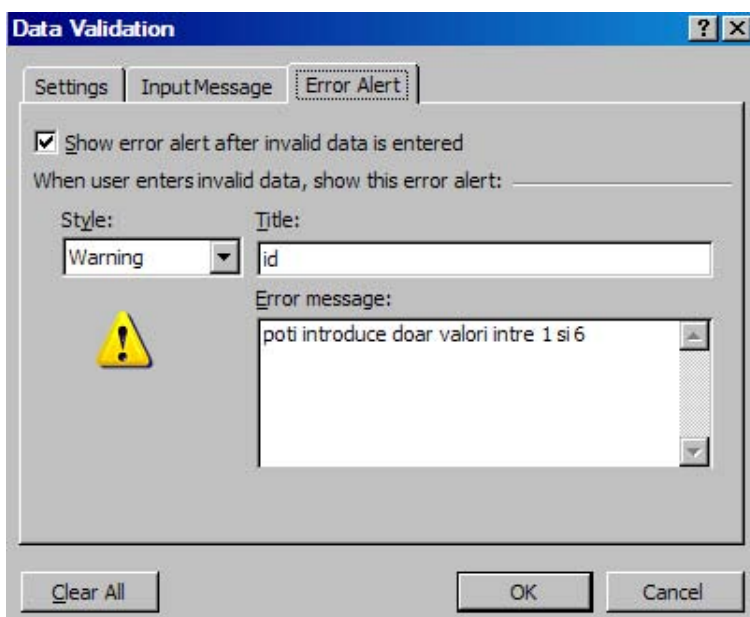
**Figura 2.3.** Impunerea condițiilor de introducere a datelor în Excel. Tabul Settings : fereastra cu condiții



În tabul **Error Alert** :

- în secțiunea **Style**, selectăm **Warning** ;
- în secțiunea **Title** tastăm numele variabilei **Id**. În această secțiune vom introduce numele variabilei pentru care impunem condiția. Astfel, vom identifica ușor la care variabilă se referă avertizarea ;
- în secțiunea **Error message** tastăm mesajul de avertizare pentru operatorul de introducere : „Poți introduce doar valori între 1 și 6” (figura 2.4). Salvăm ceea ce am lucrat (apăsăm simultan tastele CTRL + S).

**Figura 2.4.** Impunerea condițiilor de introducere a datelor în Excel. Tabul Error Alert : fereastra cu mesaj în stilul Warning



Putem trece la variabila următoare, v1. În această etapă, trebuie să instruiem programul să accepte doar codurile 1 sau 2. Vom utiliza aceleași condiții ca mai sus. La fel procedăm și cu v2. Modificăm numele variabilei în secțiunea **Title** și mesajul care apare în secțiunea **Error message**: Pentru v1, mesajul va fi „Poți introduce doar valorile 1 sau 2”. Pentru v2, mesajul va fi „Poți introduce doar valorile 1, 2 sau 3”. Salvăm ceea ce am lucrat (apăsăm simultan tastele CTRL + S).

Baza de date este finalizată. Acum trebuie să verificăm dacă funcționează conform așteptărilor. Verificarea va fi realizată prin introducerea unor coduri greșite. În acest sens, putem opta pentru două variante. Prima variantă, cea mai la îndemână, presupune să tastăm direct în celule, adică introducem valoarea 7 în celula A2. Această variantă este utilă atunci când avem puține variabile și puțini respondenți, adică sunt puține coloane și rânduri. Dacă avem foarte multe coloane și foarte multe rânduri, va deveni obositor pentru operatorul de introducere să le urmărească. Excel oferă o variantă mai simplă de folosit, constând într-un formular. Formularul poate fi accesat din meniul **Data > Form...** Înainte de a accesa acest meniu, selectăm cele trei coloane care conțin variabilele id, v1 și v2. În figura 2.5 este prezentat un formular gol și un formular cu informațiile introduse pentru respondentul 1. După ce am terminat de introdus datele pentru un respondent, apăsăm butonul **New** și trecem la următorul respondent.

**Figura 2.5.** Formular de introducere a datelor în Excel. Formular gol și formular cu informații introduse pentru un respondent

Dacă am fi introdus valoarea 7 la id, atunci programul ne-ar fi avertizat că această valoare nu face parte dintre cele valide. Repetăm operațiunea și pentru v1, introducând altceva decât valorile 1 sau 2. La fel și pentru v3, introducând altceva decât valorile 1, 2 sau 3. Dacă totul funcționează conform așteptărilor noastre, atunci putem trece la introducerea datelor. După introducere, baza de date ar trebui să arate ca în figura 2.6. Vom avea șapte rânduri pentru că primul conține numele variabilelor.

**Figura 2.6.** Bază de date creată în Excel

|   | A  | B  | C  |
|---|----|----|----|
| 1 | id | v1 | v2 |
| 2 | 1  | 1  | 1  |
| 3 | 2  | 1  | 2  |
| 4 | 3  | 1  | 2  |
| 5 | 4  | 2  | 1  |
| 6 | 5  | 2  | 3  |
| 7 | 6  | 2  | 3  |

În practică, lucrurile sunt mai complexe. De exemplu, în foarte multe chestionare, dacă nu în toate, există întrebări-filtru. Adică respondentului 1, pentru că la v2 a ales răspunsul „apă” (codul 1), ar putea să îi fie adresată o altă întrebare : „Ce marcă



preferați din lista : 1. Dorna, 2. Bucovina, 3. Izvorul minunilor ? ”. Sau responden-  
tului 2, pentru că la v2 a ales răspunsul „suc cu cofeină” (codul 2), ar putea să îi fie  
adresată o altă întrebare : „Ce marcă preferați din lista : 1. Coca Cola, 2. Pepsi  
Cola, 3. Adria Cola ? ”. În aceste situații, mi se pare mai simplu de utilizat programul  
Microsoft Access. Preferința pentru un program sau altul este, în final, o chestiune  
de gust sau de experiență cu unul sau altul. O altă caracteristică după care mă ghidez  
în alegerea programului cu care lucrez pentru o sarcină anume este ușurința cu care  
pot găsi informații ajutătoare despre diferite operațiuni pe care trebuie să le efectuez  
cu acesta. Iar Excel, Access și SPSS stau foarte bine la acest capitol.

### 2.1.2. Introducerea datelor în Microsoft Access

În cercetările la care am participat, am folosit adesea Microsoft Access, de aceea  
vă voi explica pe scurt cum se creează o bază de date în acest program.

Baza de date creată în Access sau în Excel va fi importată în SPSS. Pentru ca  
această tranziție să funcționeze corect, dar și pentru a ne fi ușor să lucrăm cu  
variabilele din baza de date, trebuie să respectăm câteva condiții :

- Să citim chestionarul cu atenție și să identificăm toate variabilele care trebuie  
să facă parte din baza de date. O întrebare poate conține mai multe variabile.  
Toate variabilele trebuie să se regăsească în baza de date.
- Chestionarele trimise în teren să aibă un identificator (id) unic. De exemplu,  
id-ul poate fi numărul chestionarului. Dacă avem 1.000 de chestionare de  
aplicat, atunci acestea sunt numerotate de la 1 la 1.000, fără repetiții. Numărul  
chestionarului va fi id-ul. Acesta poate fi și mai complex de atât, decizia pentru  
forma finală depinzând de designul cercetării. Cert este că nu există bază de  
date fără această variabilă.
- Fiecare variabilă să aibă un nume (*name*) care este diferit de al celorlalte  
variabile.
- Numele să înceapă cu o literă. Numerele pot fi folosite ulterior. Între caractere,  
fie că sunt litere, fie că sunt numere, nu se lasă spațiu. Dacă dorim să separăm  
diferite elemente ale numelui, atunci utilizăm semnul „\_”.
  - Corect : *v1*. Incorect : *1v*.
  - Corect : *v1*. Incorect : *v 1*.
  - Corect : *v\_1*. Incorect : *v 1*.
- Vă recomand să scrieți cu literă mică întregul nume. Dacă trebuie să realizați  
o analiză statistică în alt program, iar acel program face distincția între litere  
mari și litere mici, atunci există posibilitatea să vă încurcați în denumiri.
- Deși versiunile mai noi de SPSS permit să folosiți nume lungi, vă recomand  
să folosiți nume scurte, pentru a le putea găsi ușor în lista de variabile din  
meniuri. Un nume scurt este mai ușor de ținut minte decât un nume lung.

Dacă va trebui să folosiți o variabilă cu nume lung într-o analiză efectuată cu alt program de statistică, iar acel program nu acceptă decât, de pildă, maximum 8 caractere, atunci numele va fi trunchiat și s-ar putea să vă fie greu să o mai găsiți în baza de date.

Access folosește tabele și formulare create pornind de la tabele. Tabelul este baza de date. Formularul este interfața prietenoasă pe care o poate folosi operatorul pentru a introduce chestionarele în baza de date.

Prima etapă în crearea unei baze de date în Access constă în crearea unui tabel. Pentru un chestionar scurt va fi suficient un singur tabel. Pentru chestionare lungi, cu multe întrebări și, implicit, variabile, va trebui, probabil, să creați mai multe tabele. Există multe manuale și tutoriale dedicate acestui subiect. De aceea voi nota aici doar lucrurile elementare care ne interesează într-o cercetare socială obișnuită.

Să deschidem programul. Odată deschis, mergem în meniul **File > New > Blank Database**. Dăm un nume bazei de date și o salvăm undeva în computer. Inserăm un tabel în formatul **Design View**. Voi folosi ca exemplu chestionarul utilizat în cercetarea *Diagnoza calității vieții din România* (DCV 2010) realizată în 2010 de Institutul de Cercetare a Calității Vieții din cadrul Academiei Române.

Prima variabilă va fi, întotdeauna, cea care conține identificatorul unic pentru fiecare chestionar. Chestionarele au fost numerotate de la 1 la  $n$ , unde  $n$  reprezintă numărul total de chestionare completate de operatorii de teren conform designului cercetării. Această variabilă poartă numele „nrchest”. Aceasta va fi cheia primară (*primary key*) a tabelului. Putem avea o singură cheie primară într-un tabel. Access o va defini singur, dar putem să ne asigurăm că este cea corectă dacă în dreptul variabilei dorite este vizibilă o cheie. Putem alege ca această cheie să fie completată de program sau să o introducem noi. Deși a doua variantă este mai supusă greșelii, eu o prefer pentru că îmi permite să folosesc chiar informația notată pe chestionar. Acest lucru este cu atât mai important atunci când ID-ul nu pornește de la 1, ci este un cod mai complicat dat de responsabilul de teren fiecărui chestionar. Definirea manuală a cheii primare se face astfel: în tabel, în formatul **Design View**, în coloana **Field Name** introducem nrchest. În coloana **Data Type** selectăm **Number**. Am instruit, astfel, programul că pentru variabila nrchest, introducem numere. Apoi, ducem cursorul pe indicatorul rândului, dăm click dreapta și selectăm **Primary Key**. Salvăm tabelul (apăsăm simultan tastele CTRL + S).

Acum putem continua definirea variabilelor din chestionar. În chestionarul DCV 2010, respondentului îi sunt adresate mai întâi o serie de întrebări sociodemografice. Echipa care a întocmit chestionarul a avut în vedere, în faza de redactare, faptul că trebuie realizată o corespondență perfectă între hârtie și computer, între chestionar și baza de date. Astfel, toate variabilele au primit în chestionarul tipărit un nume unic care respectă condițiile enumerate mai sus.

Prima variabilă are numele d1, a doua d2, a treia d3, iar lista continuă până la d119. După d119, urmează o secțiune scurtă de întrebări adresate operatorului de teren, acestea având numele op1, op2, ..., op9.

Primul lucru pe care respondentul este rugat să îl declare este genul. Numele acestei variabile este d1. Genul (d1) are două variante de răspuns : masculin sau feminin. Varianta masculin a primit codul 1. Varianta feminin a primit codul 2. În baza de date trebuie introduse codurile care se regăsesc în chestionar și nimic altceva. Transformările se fac, ulterior, în SPSS. De exemplu, dacă în chestionar respondentului i s-a cerut să declare anul în care s-a născut, atunci în baza de date vom introduce anul nașterii. Nu îi vom cere operatorului de introducere să calculeze vârsta și să introducă valoarea rezultată.

Revenind în Access, în tabelul în format **Design View**, pe următorul rând, sub nrchest, vom introduce în coloana **Field Name** d1, iar în coloana **Data Type** selectăm **Number**. Revin la modul de formatare a chestionarului. Chestionarul este folosit pentru că vrem să calculăm anumite statistici. Statisticile pe care vrem să le calculăm constituie o decizie pe care, teoretic, cercetătorul o ia înainte de a trimite chestionarul în teren. Astfel, vă asigurați că se vor culege informațiile de care aveți nevoie pentru a răspunde la întrebarea de cercetare. Pentru că statisticile se calculează folosind numere, atunci, în chestionar, când folosiți aplicarea față în față cu un operator de teren, din punctul meu de vedere, este obligatoriu să notați codurile atribuite variantelor de răspuns :

| Corect  | Inc corect   |
|---|--|
| D1. Sexul :<br>1. masculin<br>2. feminin  | 1. Sexul :<br><input type="checkbox"/> masculin<br><input type="checkbox"/> feminin  |
| D4. Statutul ocupațional :<br>1. salariat<br>2. pe cont propriu<br>3. patron<br>4. zilier   | 4. Statutul ocupațional :<br><input type="checkbox"/> salariat<br><input type="checkbox"/> pe cont propriu<br><input type="checkbox"/> patron<br><input type="checkbox"/> zilier   |
| D26. Cum caracterizați calitatea transportului în comun în localitatea dvs :<br>1. foarte proastă<br>2. proastă<br>3. satisfăcătoare<br>4. bună<br>5. foarte bună | 26. Cum caracterizați calitatea transportului în comun în localitatea dvs :<br><input type="checkbox"/> foarte proastă<br><input type="checkbox"/> proastă<br><input type="checkbox"/> satisfăcătoare<br><input type="checkbox"/> bună<br><input type="checkbox"/> foarte bună |

*Sursa : chestionarul **Diagnoza calității vieții în România 2010**, ICCV.*

Dacă nu notăm codurile, atunci operatorul de introducere a datelor va trebui fie să alocă mult timp înainte de a trece la introducerea efectivă, pentru notarea pe chestionare a codurilor aferente fiecărei variante de răspuns, fie să fie atent ca la fiecare variabilă să introducă corect codul. Se pierde, astfel, timp prețios

și crește riscul apariției erorilor de introducere. O altă eroare de redactare observată în acest exemplu este atribuirea numelor exclusiv sub formă de număr : în loc de D1, D4 sau D26 (am păstrat numerotarea din chestionar), cercetătorul a atribuit doar 1, 4 sau 26.

Tipul de variabilă (**Data Type**) depinde de caracteristicile informațiilor conținute. De regulă, introducem numere și, uneori, text. Pentru fiecare variabilă definită în **Design View** trebuie să alegem un set de proprietăți, dintre care le prezint pe cele mai importante :

- **Field Size.** De regulă, vom alege între **Byte**, **Integer** sau **Long Integer**. Diferența dintre ele constă în numărul de cifre pe care le poate avea valoarea introdusă.
- **Default value.** Dacă operatorii de teren ar lucra perfect la aplicarea chestionarelor, atunci toate celulele din baza de date ar avea informații conform instrucțiunilor chestionarului. Adică ar fi introduse fie răspunsurile valide, fie codurile pentru nonrăspuns. Nonrăspunsul este de trei tipuri : respondentul refuză să răspundă, respondentul nu știe să răspundă sau întrebarea nu trebuie să îi fie aplicată respondentului. Acestea primesc coduri speciale, diferite semnificativ ca formă de codurile valide. Cele mai utilizate în România sunt 97 = Nu este cazul (NC), 98 = Nu știu (NS), 99 = Nu răspund (NR). Există situații în care trebuie să le transformăm. De exemplu, o femeie nu vrea să își declare vârsta. Operatorul de teren ar trebui să noteze pe chestionar codul 99. Dar 99 poate fi o vârstă validă. Atunci, echipa de cercetare, sub îndrumarea celui care face designul bazei de date, ar putea să instruiască operatorul să noteze pe chestionar codul 999. Aceasta nu mai este o vârstă validă. Dar, dacă ne gândim la salariul lunar, 999 lei poate fi un salariu valid. Atunci, codul de nonrăspuns ar putea deveni -1. Acesta nu mai este un salariu valid. Ideea este să folosim un cod cu totul diferit de variantele de răspuns valide. Punând unul dintre aceste coduri ca **Default Value**, îi spunem programului să introducă singur valoarea respectivă. În acest mod, ne asigurăm că am definit un răspuns ușor de înțeles, când începem analiza statistică. Dacă lăsăm celula goală în tabel, în această fază, s-ar putea să nu mai știm ce am vrut de fapt să simbolizeze : este o lipsă de răspuns, este o scăpare a operatorului de introducere etc. ? Alegerea codului pentru **Default Value** depinde de tipul întrebării. Dacă răspunsul la întrebare nu depinde de un filtru, atunci vom folosi codul 99 (NR). Dacă răspunsul la întrebare depinde de un filtru, atunci vom folosi codul 97 (NC).
- **Validation Rule.** În acest câmp, introducem o condiție prin care instruiem programul să accepte doar codurile valide înregistrate în chestionar. De exemplu, la d1 avem trei coduri valide : 1 = masculin, 2 = feminin și 99 = nu răspund. În practică, ultimul cod nu este acceptabil, pentru că operatorul trebuie să vină cu informații complete măcar la variabilele sociodemografice esențiale. Așadar, regula noastră de validare va fi „1 Or 2 Or 99”. Practic, îi spunem programului să primească doar codurile 1, 2 sau 99. Dacă introducem

codul 3, nu ne va permite să mergem mai departe, deci trebuie să corectăm regula de validare. În chestionarul DCV 2010, la variabila D26, avem cinci variante de răspuns care au primit coduri de la 1 la 5. Am putea scrie „1 Or 2 Or 3 or 4 Or 5 or 99”. Dar, mai simplu, putem scrie : „Between 1 and 5 Or 99”.

- **Validation Text.** Aici putem, opțional, să punem un mesaj ajutător pentru operatorul de introducere. De exemplu : „Variantele corecte sunt 1, 2 sau 99”. Operatorul va identifica mai repede eroarea pe care a realizat-o.

Să recapitulăm luând un exemplu care include și un filtru în chestionar. După precizarea genului, în DCV 2010, respondentul este rugat să declare care este ocupația sa principală actuală. Această variabilă are numele d2, zece variante de răspuns, fiecare având un cod unic și două tipuri de filtre :

|                                   |   |                                |
|-----------------------------------|---|--------------------------------|
| 1. agricultor                     |   |                                |
| 2. muncitor (meserias)            |   |                                |
| 3. tehnician, maistru, funcționar |   |                                |
| 4. ocupație cu studii superioare  |   |                                |
| 5. altă ocupație                  | → | Care ? _____                   |
| 6. elev, student                  | → | Dacă 6, sari la întrebarea d6. |
| 7. pensionar                      | → | Dacă 6, sari la întrebarea d6. |
| 8. casnică                        | → | Dacă 6, sari la întrebarea d6. |
| 9. acum sunt șomer                | → | Dacă 6, sari la întrebarea d6. |
| 10. patron                        |   |                                |

Ordinea variabilelor din tabel trebuie să respecte ordinea variabilelor din baza de date. Așadar, următorul rând în **Design View**, după d1, va deveni d2. La **Default Value** vom introduce 99. La **Validation Rule** vom scrie „Between 1 And 10 Or 99”. La **Validation Text** vom scrie „Poți introduce doar coduri între 1 și 10 sau 99”.

Dacă respondentul alege una dintre variantele 1, 2, 3, 4 sau 10, atunci i se va adresa întrebarea următoare : d3. „Din ce an aveți această ocupație ?”. Dacă respondentul alege varianta 5, atunci va trebui să completeze răspunsul la întrebarea „Care ?”. Dacă respondentul alege una dintre variantele 6, 7, 8 sau 9, atunci întrebările d3, d4 și d5 nu i se aplică și se trece direct la întrebarea d6. Cele trei întrebări nu i se aplică, pentru că se referă la ocupație. Așadar, avem mai multe filtre care trebuie definite și în baza de date. Folosim filtre în baza de date pentru a grăbi procesul introducerii : introducând valoarea automată 97 la variabilele corespunzătoare, putem sări peste acestea, scutind timp pe care îl putem alocă analizei statistice propriu-zise.

Nu definim filtrele în tabel, ci în formular. Formularul va fi elaborat după ce a fost finalizat tabelul. Adică definim toate variabilele și proprietățile lor în tabel, salvăm și abia apoi trecem la formular. Acum ardem etapele doar în scop didactic.

Să presupunem, așadar, că am finalizat tabelul introducând toate variabilele din chestionar. Formularul este inserat și deschis tot în formatul **Design View**. În principiu, toți acești pași sunt intuitivi în interfața programului, motiv pentru care nu mai insist aici. Am creat formularul care conține toate variabilele din

tabel. Am putea să îi aducem tot felul de îmbunătățiri estetice, dar aceasta este o chestiune de gust, și nu de o necesitate imperioasă. Introducerea filtrelor este însă foarte importantă.

Primele două filtre se stabilesc pentru variabila d2 :

- dacă respondentul răspunde cu 5 la d2, atunci programul trebuie să meargă la variabila „Care ?”.
- dacă respondentul răspunde cu 6, 7, 8 sau 9, atunci programul trebuie să sară peste calupul de întrebări dintre d2 și d6, mergând direct la d6.

În chestionar, variabila „Care ?” nu a primit un nume cum ar fi d1, d2, d3 etc. Această situație poate fi remediată ușor în program introducând în tabel numele „d2care”. Pentru că respondentului i s-a cerut să precizeze cu propriile lui cuvinte ce ocupație are, răspunsurile sunt înregistrate sub formă de text. În câmpul **Data Type** alegem fie **Text**, fie **Memo**. Opțiunea între **Text** și **Memo** depinde de numărul de caractere care va fi introdus. Pentru simplitate, eu prefer să le definesc pe toate **Memo**. La **Default Value** am introdus 97 pentru că această întrebare se aplică doar celor care au ales codul 5 la d2. Pentru cei care au ales codurile 1-4, respectiv 6-10, această întrebare nu se aplică.

În formular, intrăm în modul **Design View**. Mergem la d2 și dăm click dreapta pe celulă (nu pe etichetă). Alegem opțiunea **Properties**. Se va deschide o fereastră din care, pentru această situație, ne interesează tabul **Event**. Din tabul **Event** ne interesează rândul **On Exit**. Practic, acest eveniment instruește programul să aleagă o acțiune în funcție de codul introdus în d2 atunci când apăsăm tasta **Tab** sau tasta **Enter**, adică trecem la următoarea variabilă din bază. În rândul **On Exit** selectăm **Event Procedure**, apoi dăm click pe cele trei puncte din dreapta celulei. Se deschide o fereastră de cod. Între **Private sub...** și **End sub**, trebuie să introducem sintaxa :

```
If Me![d2] = 1 Then
Me![d3].SetFocus
ElseIf Me![d2] = 2 Then
Me![d3].SetFocus
ElseIf Me![d2] = 3 Then
Me![d3].SetFocus
ElseIf Me![d2] = 4 Then
Me![d3].SetFocus
ElseIf Me![d2] = 5 Then
Me![d2care].SetFocus
ElseIf Me![d2] = 6 Then
Me![d6].SetFocus
ElseIf Me![d2] = 7 Then
Me![d6].SetFocus
ElseIf Me![d2] = 8 Then
```

```

Me ! [d6].SetFocus
ElseIf Me ! [d2] = 9 Then
Me ! [d6].SetFocus
ElseIf Me ! [d2] = 10 Then
Me ! [d3].SetFocus
ElseIf Me ! [d2] = 99 Then
Me ! [d6].SetFocus
End If

```

Această sintaxă instruește programul să respecte filtrele :

- dacă la d2 primește codurile 1-4 sau 10, să treacă la variabila d3 pentru că aceasta se aplică acestor respondenți ;
- dacă la d2 primește codul 5, să treacă la variabila d2care, pentru că aceasta se aplică acestor respondenți ;
- dacă la d2 primește codurile 6-9 sau 99, să meargă la variabila d6, pentru că aceasta se aplică acestor respondenți. Trebuie să definim condiția și pentru codul de nonrăspuns.

Salvăm și închidem fereastra de cod. Ne întoarcem în formular (nu uităm că tabelul este finalizat deja și nu mai intervenim în el, decât în situații excepționale) și continuăm cu celelalte variabile, dacă este cazul.

În formular, în modul **Design View**, putem modifica și estetica formularului. Putem introduce etichete pentru calupuri de întrebări, săgeți ajutătoare pentru operator etc. Mai important mi se pare să avem în vedere că responsabilul cu elaborarea bazei de date poate lucra cu o versiune mai nouă/veche a programului, iar operatorii de introducere cu una mai veche/nouă a acestuia. Acesta trebuie să asigure compatibilitatea între versiuni. De preferat ar fi să se lucreze pe aceeași versiune.

Din punct de vedere estetic, mi se pare important ca formularul să aibă variabilele dispuse în așa fel încât să încapă pe o jumătate de ecran. Prefer ca introducerea să decurgă de sus în jos, adică variabilele să fie una sub alta :

| Varianta preferată de mine |        |   |     | Variantă posibilă |        |     |     |
|----------------------------|--------|---|-----|-------------------|--------|-----|-----|
| ↓                          | d1     |   | d3  | è                 | d1     | ↙   | d2  |
|                            | d2     |   | d4  |                   | d2care | ↙   | d3  |
|                            | d2care |   | ... |                   | d4     | ↙   | d5  |
|                            | ...    | ↗ | ... |                   | ...    | ... | ... |

După ce am realizat formularul și am introdus toate condițiile, trebuie să verificăm dacă am lucrat corect. Acest lucru se face simplu, după cum am discutat și la Excel, introducând în celulele formularului valori ce nu se regăsesc printre răspunsurile valide sau care nu sunt coduri de nonrăspuns. Vom observa imediat dacă filtrele funcționează sau nu.



## 2.2. Cum ducem datele în SPSS ?

Am încheiat introducerea datelor. Trebuie să trecem la etapa de curățare a bazei de date, pe care o realizăm în SPSS. Așadar, trebuie să ducem datele din formatul Excel sau Access în formatul SPSS. Există mai multe posibilități în acest sens, dar, pentru că „paza bună trece primejdia rea”, prefer ca, mai întâi, să vizualizez datele în Excel, iar din Excel să le duc în SPSS.

Din Access ducem datele în Excel astfel :

- selectăm tabelul pe care dorim să îl exportăm în Excel,
- deschidem meniul **File > Export**,
- în fereastra care se deschide, la secțiunea **Save as type** alegem unul dintre formatele Excel, de exemplu, **Microsoft Excel 97-2002**, dacă lucrăm cu versiunea 2002 a Access,
- denumim tabelul în modul dorit și apăsăm tasta **Enter** sau butonul **Export**.

Acum datele sunt în formatul Excel. Pentru a evita erorile generate de modul cum tratează SPSS informația venită din alte programe, mai ales în versiunile mai vechi, recomand să verificați dacă :

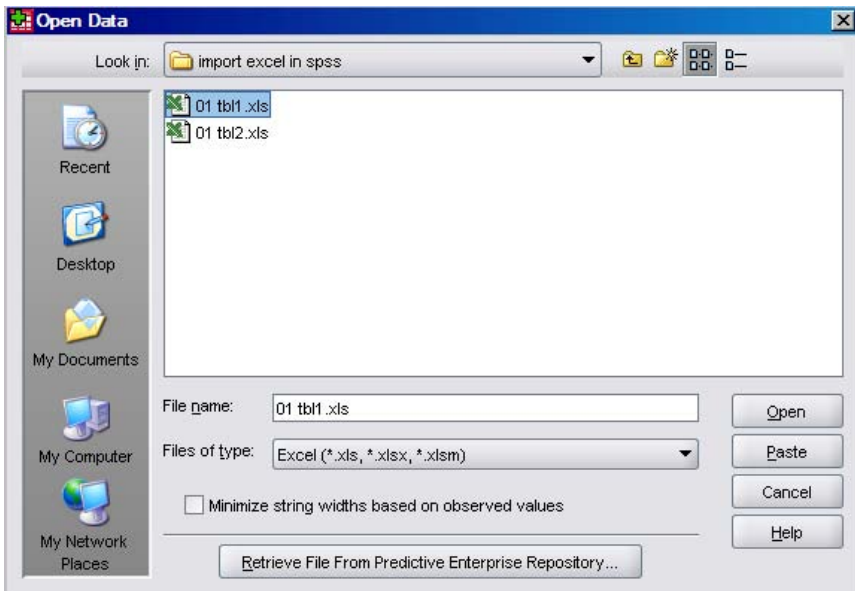
- primul rând din Excel conține numele variabilelor, iar numele respectă condițiile enunțate în acest capitol ;
- toate numerele din celule sunt tratate de Excel ca fiind numere. Pentru siguranță, putem selecta variabilele care conțin numere (coloanele din Excel) dând click dreapta pe numele coloanei și selectând **Format Cells**. În fereastra care se deschide, în tabul **Number**, la secțiunea **Category** alegem **Number**. De regulă, mai ales în versiunile noi, SPSS citește corect informația din Excel, aceasta fiind o măsură de precauție.

Suntem pregătiți să ducem datele în SPSS. Pentru exemplificare, voi folosi date din DCV 2010. În SPSS pot fi importate baze de date salvate și în alte formate (de exemplu, fișierele care au extensia **.csv** sunt deseori folosite de analiști) sau chiar în formatul specific altui program de statistică cum ar Stata (fișiere cu extensia **.dta**). SPSS citește și aceste extensii. Lista completă a formatelor recunoscute de SPSS poate fi găsită în documentația programului. Mai poate fi utilizat un program comercial, Stat Transfer, care este dedicat acestui gen de operațiuni. Flexibilitatea este destul de ridicată în domeniul programelor statistice, odată ce ajungi să te familiarizezi cu limbajul acestora.

Deschidem programul și mergem în meniul **File > Open > Data**. În fereastra care se deschide, selectăm locul unde am salvat tabelele în format Excel. Apoi, în secțiunea **Files of type**, alegem **Excel (\*.xls, .xlsx, .xlsm)**. Inițial, este selectat **SPSS Statistics (\*.sav)**. Selectăm tabelul Excel pe care vrem să îl importăm în SPSS (figura 2.7).



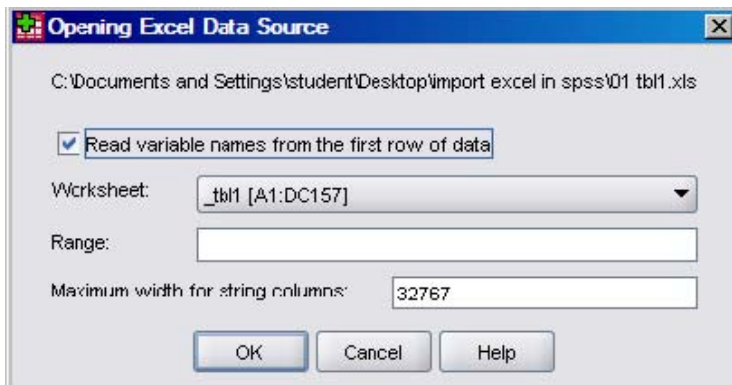
**Figura 2.7.** Importarea datelor din formatul Excel în formatul SPSS.  
Selectarea fișierului Excel



După ce apăsăm butonul **Open**, se va deschide fereastra din figura 2.8. Aici trebuie să selectăm foaia de lucru în care sunt datele care ne interesează. De regulă, avem o singură foaie de lucru. Dar dacă avem mai multe și ne interesează una anume, o vom selecta din listă pe cea corespunzătoare. Înainte de a apăsa OK, verificați dacă celula **Read variable names from the first row of data** este selectată. Ar trebui să fie.

Salvăm baza de date rezultată (apăsăm simultan tastele CTRL + S).

**Figura 2.8.** Importarea datelor din formatul Excel în formatul SPSS.  
Selectarea foii de lucru care conține datele



Repetăm acești pași ori de câte ori este nevoie. De exemplu, la DCV 2010 au introdus chestionare mai mulți operatori. Întrucât chestionarul utilizat are multe variabile, dată fiind complexitatea temei, a fost nevoie de două tabele în Access care cuprindeau, separat, aproximativ jumătate din chestionar. Așadar, avem de importat în SPSS două baze de date în Excel de la fiecare operator pentru fiecare dintre cei șase operatori. Rezultatul final al procesului de importare trebuie să fie o singură bază de date în SPSS. De aceea, trebuie să parcurgem următoarele etape :

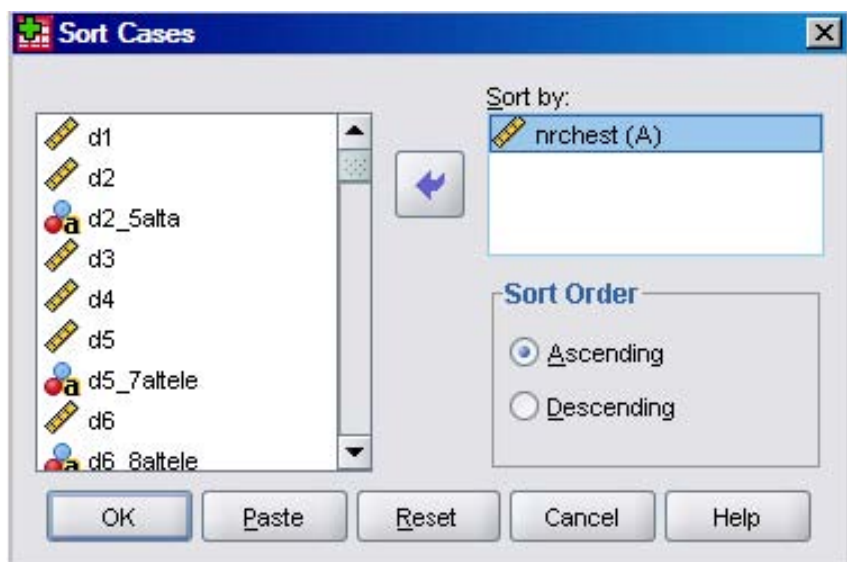
- urmând pașii descriși anterior, importăm din Excel în SPSS, pentru fiecare operator, cele două baze de date : una cu prima parte a chestionarului și una cu a doua parte a chestionarului ;
- în cazul fiecărui operator, vom uni în SPSS cele două baze pentru a avea o singură bază, adică întregul chestionar. Rezultă, astfel, șase baze în SPSS pentru toți operatorii ;
- în fine, vom uni în SPSS cele șase baze într-una singură. Aceasta este baza pe care vom realiza procesul de curățare.

Să le luăm pe rând. Pentru operatorul M, avem două baze în SPSS. La fel pentru operatorul A, D, E etc. Mai întâi, vom lucra cu cele două baze ale operatorului M. Aceste baze au aceiași respondenți, dar variabile diferite. Am afirmat mai sus că a fost nevoie să împărțim chestionarul în Access, dat fiind numărul mare al variabilelor din chestionarul complex. Pentru a uni aceste două baze, folosim meniul **Data > Merge Files > Add Variables**. Obligativ, ambele baze vor avea o variabilă de identificare care ia valori unice pentru fiecare respondent. Fără ea, unirea nu se poate face corect. Aici, această variabilă este numărul chestionarului care, în ambele baze de date, poartă numele nrchest. Deschidem ambele baze de date în SPSS. Primul lucru pe care îl facem este să ordonăm bazele de date, în aceeași direcție, crescător, după nrchest : **Data > Sort Cases >** trecem nrchest în dreapta folosind săgeata > lăsăm bifat **Ascending > OK** (figura 2.9). Mai rapid, putem să deschidem **Data View**, dăm click dreapta pe numele variabilei nrchest și selectăm **Sort Ascending**. Salvăm ambele baze de date (apăsăm simultan tastele CTRL + S).

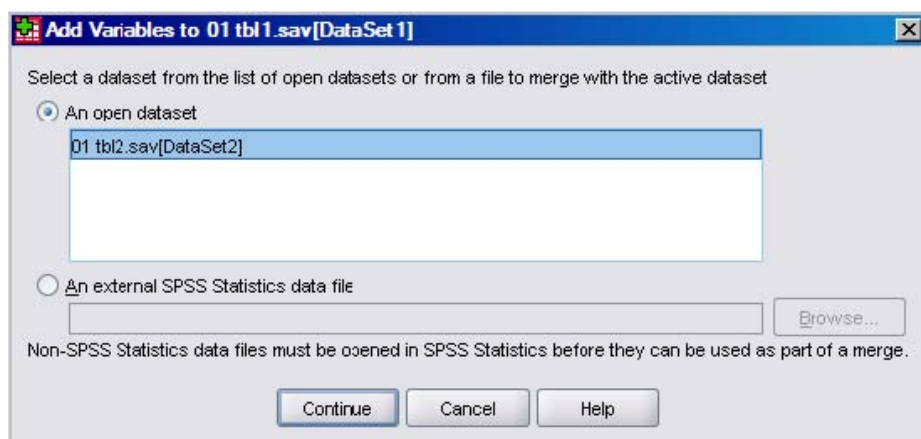
Acum putem ține deschise ambele baze de date sau doar pe cea în care aducem noile variabile. Voi explica cum procedăm pentru prima variantă. Deschidem baza primară în meniul **Data > Merge Files > Add Variables**. Se deschide fereastra din figura 2.10. Selectăm cea de-a doua bază. Dacă această bază nu era deschisă, trebuia să selectăm **An external SPSS Statistics data file** și să căutăm pe computer unde este salvată. Apăsăm **Continue**. În fereastra care se deschide, bifăm **Match cases on key variables in sorted files**, trecem nrchest în căsuța **Key Variables** și apăsăm **OK**. Programul ne va avertiza că

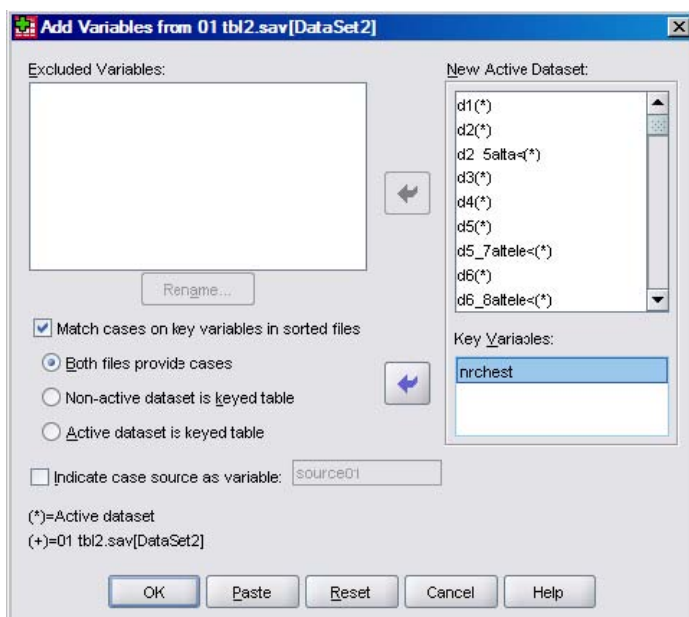
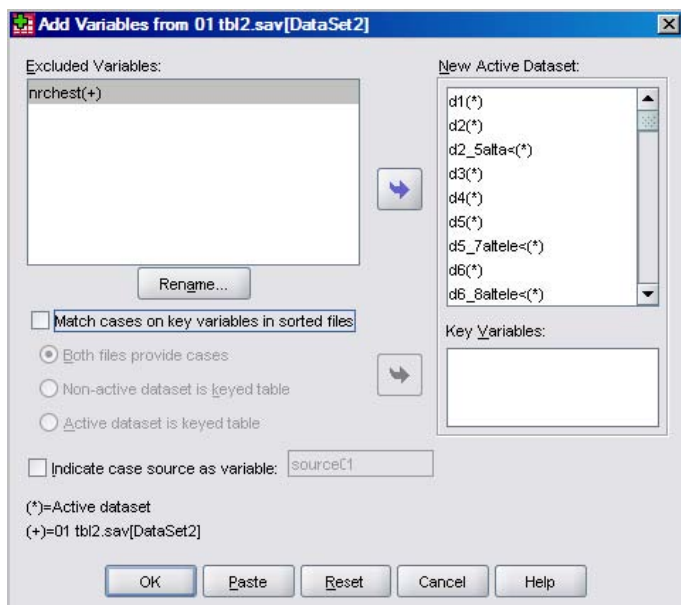
trebuia să ordonăm ambele baze de date după variabila-cheie, aici nrchest. Din acest moment, datele din cele două baze se află într-una singură. Dacă dorim să păstrăm bazele inițiale și să avem separat baza unită, atunci va trebui să salvăm rezultatul sub o altă denumire.

**Figura 2.9.** Ordonarea cazurilor : după o variabilă, de la valorile mici la valorile mari



**Figura 2.10.** Unirea a două baze cu aceiași respondenți și variabile diferite





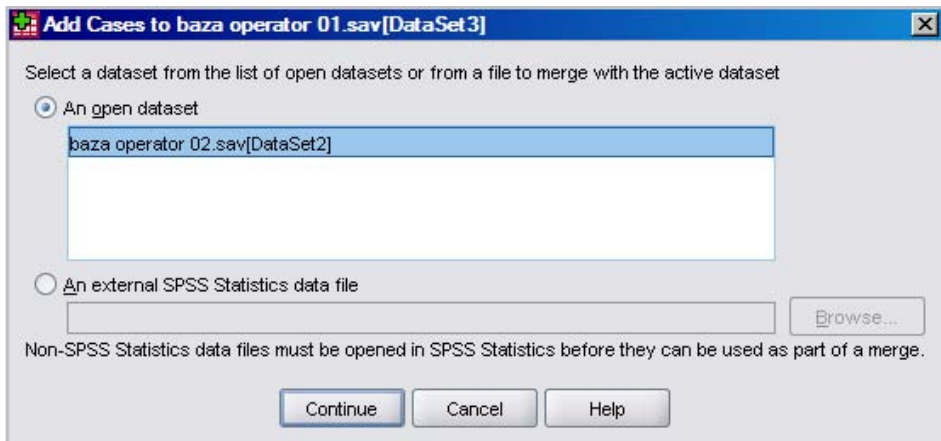
Realizăm această operație pentru toți cei șase operatori de introducere.

După ce am încheiat activitatea, trebuie să unim cele șase baze de date rezultate. De data aceasta, variabilele sunt aceleași, însă diferă respondenții. Vom uni loturi de respondenți sau, altfel spus, de chestionare. Vom utiliza meniul **Data > Merge Files > Add cases**.

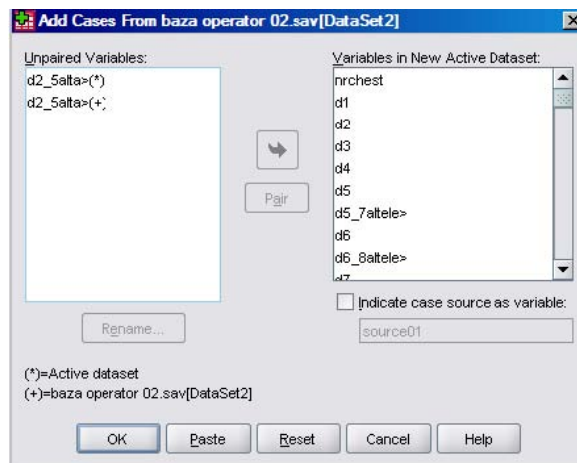
Deschidem baza de date primară, cea în care adăugăm restul de cazuri. În cazul nostru, am denumit operatorii folosind numere, pentru a nu ne încurca în denumiri : baza operator 01.sav, baza operator 02.sav etc. Voi utiliza, ca bază primară, baza primului operator. O deschidem și o ordonăm ascendent cazurile în funcție de variabila de identificare, nrchest. Deschidem baza de date a următorului operator. Ordonăm ascendent. Salvăm ambele baze, după această operație (apăsăm simultan tastele CTRL + S). Apoi revenim la baza primară, fără să o închidem pe cealaltă. Mergem în meniul **Data > Merge Files > Add cases**. Se deschide fereastra din figura 2.11a. Selectăm baza pe care dorim să o adăugăm în baza primară.

**Figura 2.11.** Unirea a două baze cu respondenți diferiți și aceleași variabile

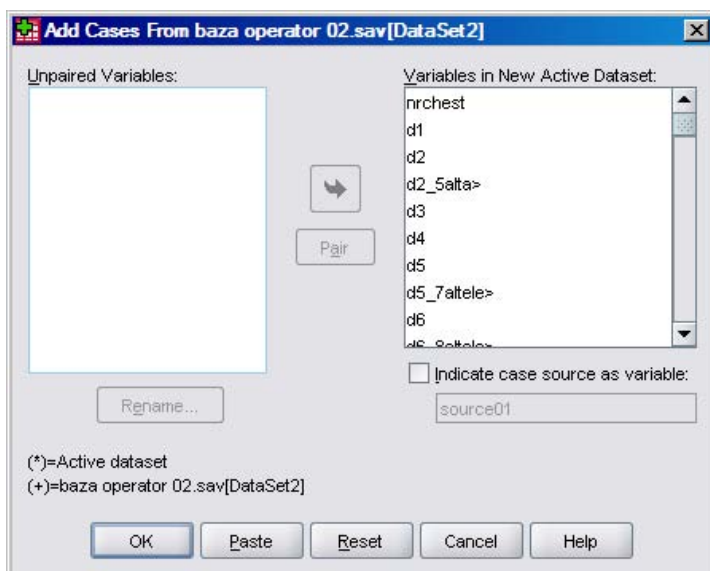
(a)



(b)



(c)



Dacă în secțiunea **Unpaired Variables** (figura 2.11b) apar variabile, înseamnă că nu putem continua unirea. Variabilele care apar în această fereastră au proprietăți diferite în cele două baze care trebuie unite. Proprietățile variabilelor pot fi vizualizate în **Variable View**. Renunțăm momentan la unire și căutăm sursa problemei. În **Variable View**, pentru d2\_5alta, în baza operator 01.sav, în coloana **Width** observăm valoarea 765, iar în baza operator 02.sav, observăm valoarea 1000 (figura 2.12). Aceasta este o variabilă care conține text, mai exact, răspunsul la întrebarea „Care ocupație?” adresată celor care au ales răspunsul „Altă ocupație” la d2, „Ocupația dvs principală (actuală)”. Situația se poate repeta pentru toți operatorii de introducere a datelor. De aceea, înainte de a încerca să unim bazele de date, trebuie să ne uităm în fiecare, în **Variable View** la ce scrie în dreptul variabilelor desperecheate (*unpaired*). Alegem o valoare comună pentru toate cele șase baze, și anume pe cea mai mare. Operăm modificările și revenim în meniul de unire. Dacă am lucrat corect, atunci celula **Unpaired Variables** ar trebui să fie goală (figura 2.11c). Apăsăm OK. Repetăm operația, până unim toate cele șase baze de la cei șase operatori.

**Figura 2.12.** Variabile cu proprietăți diferite (coloana Width din Variable View)

|   | Name     | Type    | Width |
|---|----------|---------|-------|
| 1 | nrchest  | Numeric | 8     |
| 2 | d1       | Numeric | 8     |
| 3 | d2       | Numeric | 8     |
| 4 | d2_5alta | String  | 765   |

|   | Name     | Type    | Width |
|---|----------|---------|-------|
| 1 | nrchest  | Numeric | 8     |
| 2 | d1       | Numeric | 8     |
| 3 | d2       | Numeric | 8     |
| 4 | d2_5alta | String  | 1,000 |

Din acest moment, putem începe operațiunea de curățare și de validare a bazei de date, acesta fiind subiectul capitolului 4.

## 2.3. Exerciții

Pentru aceste exerciții, utilizăm baza de date și/sau chestionarul World Values Survey 2012 rezultate în urma aplicării chestionarului în România. Baza de date poate fi descărcată de pe pagina de internet a Grupului Românesc pentru Studiul Valorilor Sociale (<http://www.romanianvalues.ro>).

1. Deschideți chestionarul WVS 2012. Alegeți, la întâmplare, două pagini din chestionar. Răspundeți la întrebările de pe aceste două pagini.
2. Realizați în Excel o bază de date care să corespundă acestor două pagini de chestionar.
3. Introduceți răspunsurile dvs. în baza de date pe care ați creat-o.
4. Importați baza de date în SPSS.
5. Rugați un coleg să vă răspundă la cele două pagini de întrebări selectate anterior.
6. Introduceți răspunsurile colegului într-o bază de date diferită de cea în care se află răspunsurile dvs.
7. Importați baza de date cu răspunsurile colegului în SPSS.
8. Uniți cele două baze de date.
9. Alegeți la întâmplare alte două pagini din chestionar. Răspundeți la întrebările de pe aceste două pagini. Rugați același coleg să vă răspundă și la aceste întrebări.
10. Realizați în Excel o bază de date care să corespundă acestor două pagini din chestionar.
11. Introduceți în baza de date creată răspunsurile dvs. și ale colegului dvs.
12. Importați baza de date în SPSS.
13. Uniți această bază de date cu cea obținută anterior în SPSS.





### 3. Gestionarea bazei de date

Manipularea și gestionarea bazei de date presupun un set de cunoștințe indispensabile analistului. Acesta trebuie să știe cum se ponderează (*weighting*) o bază de date, cum se filtrează (*select cases*) sau le separă (*split file*), cum se agregă (*aggregate*) sau se restructurează (*restructure*) etc.

Pentru începători, cred că cele mai importante operațiuni sunt cele de ponderare, filtrare și separare. Meniurile aferente acestora și pe care le discut în acest capitol sunt : **Data > Weight Cases**, **Data > Select Cases** și **Data > Split File**.

Ponderarea se referă la ajustarea bazei de date astfel încât structura eșantionului pentru variabile-cheie să fie similară cu structura populației din care a fost extras acesta și pentru care dorim să facem inferențe. Filtrarea este folosită atunci când dorim să lucrăm doar cu anumite cazuri din baza de date sau să realizăm o nouă bază de date, mai restrânsă decât cea inițială. Separarea este folosită atunci când dorim să rulăm o analiză pentru grupuri diferite și să comparăm rezultatele într-o singură fereastră.

Mai întâi vom prezenta câteva setări ale programului care ni se par utile pentru că vă ajută să vizualizați mai bine informația conținută în baza de date atunci când consultați **Outputul**.

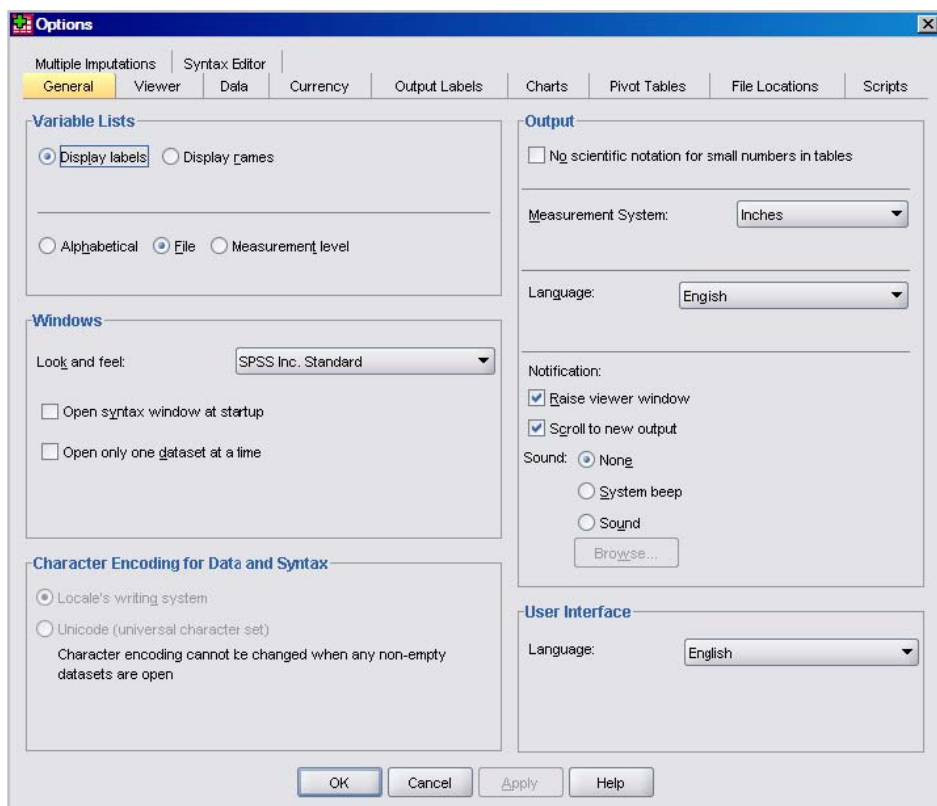
#### 3.1. Câteva setări elementare (Edit > Options)

SPSS este apreciat, printre altele, pentru că are o interfață simplă care îi permite utilizatorului să găsească rapid lucrurile de care are nevoie. În această secțiune, prezint câteva setări care cresc ușurința cu care se poate utiliza interfața. Aceste setări pot fi accesate și modificate în meniul **Edit > Options**. Figura 3.1a prezintă fereastra care apare când deschidem acest meniu. Fiind un program complex, și opțiunile sunt numeroase. Aspectul pozitiv este că ne sunt permise destul de multe intervenții în opțiunile programului, astfel încât să îl putem ajusta conform nevoilor și preferințelor noastre. Cele la care mă opresc sunt preferințele la care am ajuns în timp utilizând programul. Ați putea avea și altele pe măsură ce dobândiți experiență cu programul.

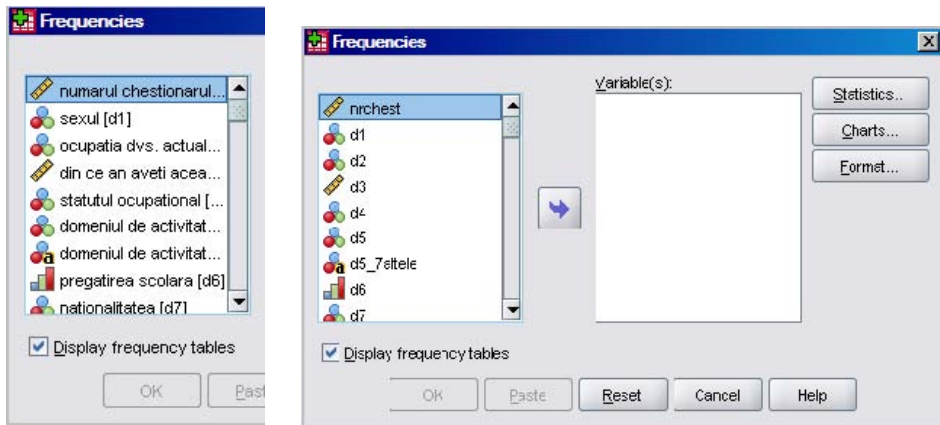
Dintre taburile de opțiuni, ne interesează următoarele : **General**, **Output Labels**, **Pivot Tables**, **File Locations** și **Syntax Editor**.

În tabul **General**, la secțiunea **Variable Lists**, bifăm **Display names**. Inițial, este bifat **Display labels**. Această operațiune va permite ca, atunci când deschidem meniurile de analiză, să observăm în lista de variabile numele în locul etichetei (*label*). Putem observa diferența în figura 3.1b care prezintă meniul **Analyze > Descriptive Statistics > Frequencies**. În ceea ce mă privește, când deschid meniul pentru analize, atunci când văd numele, nu eticheta, îmi este mult mai ușor să găsesc variabilele în lista de variabile. De altfel, putem căuta rapid, după nume, orice variabilă : dăm click în lista de variabile (coloana din stânga) și tastăm rapid primele două-trei caractere din numele acesteia. În versiunile mai noi de SPSS, putem trece foarte ușor, chiar în interiorul ferestrelor de analize, între nume și etichete. Mergem în lista de variabile, dăm click dreapta pe oricare variabilă și alegem **Display Variable Names** sau **Display Variable Labels**, în funcție de preferințe și nevoi (figura 3.1c).

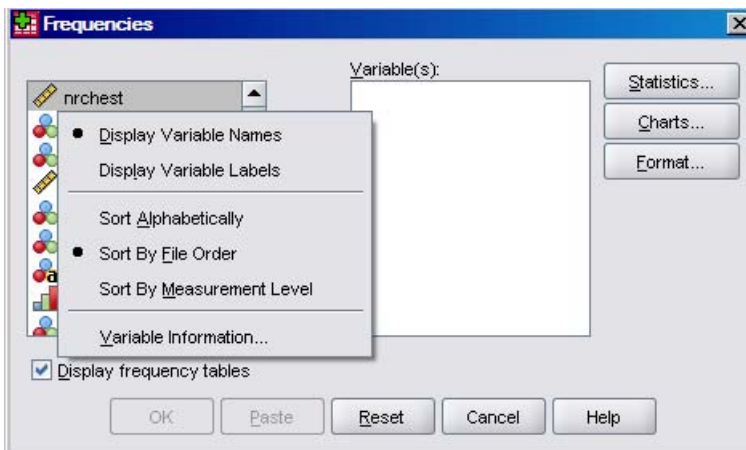
**Figura 3.1.** Setări care cresc ușurința de utilizare a programului. Tabul General (a)



(b)



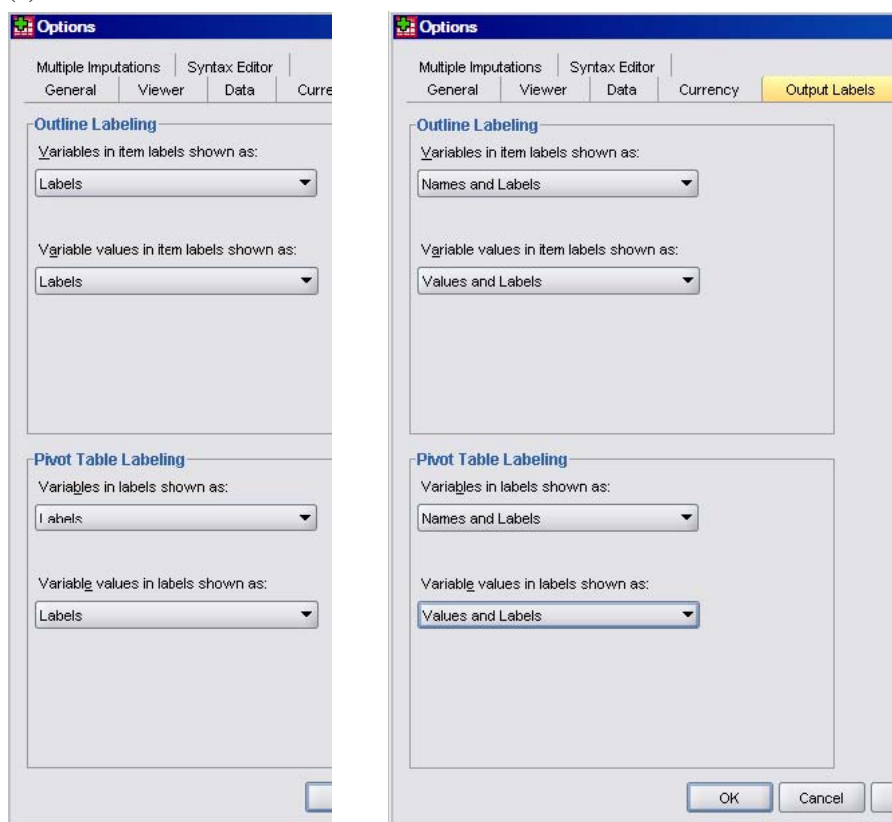
(c)



În tabul **Output Labels** (figura 3.2a) vom selecta, pentru fiecare dintre cele patru câmpuri din secțiunile **Outline Labeling** și **Pivot Table Labeling**, ambele variante: **Names and Labels**, respectiv **Values and Labels**. Făcând acest lucru, în **Output** vor fi afișate, simultan, atât numele, cât și eticheta variabilei, respectiv codurile și etichetele codurilor atribuite variantelor de răspuns. În figura 3.2b este prezentat rezultatul ambelor opțiuni. În al doilea tabel, după ce am modificat opțiunile respective, observăm atât numele, cât și eticheta variabilei, odată cu codurile și etichetele atribuite acestora.

**Figura 3.2.** Tabul Output Labels (Edit > Options) : două tipuri de vizualizare în Output

(a)



(b)

**ocupatia dvs. actuala (principala)**

|                                | Frequency | Percent | Valid Percent | Cumulative Percent |
|--------------------------------|-----------|---------|---------------|--------------------|
| Valid agricultor               | 259       | 22.3    | 22.3          | 22.3               |
| muncitori (meserias)           | 247       | 21.3    | 21.3          | 43.6               |
| tehnician, maistru, functionar | 74        | 6.4     | 6.4           | 50.0               |
| ocupatii cu studii superioare  | 106       | 9.1     | 9.1           | 59.1               |
| elev, student                  | 52        | 4.5     | 4.5           | 63.6               |
| pensionar                      | 267       | 23.0    | 23.0          | 86.6               |
| casnica                        | 62        | 5.3     | 5.3           | 91.9               |
| acum sunt somer                | 82        | 7.1     | 7.1           | 99.0               |
| patron                         | 12        | 1.0     | 1.0           | 100.0              |
| Total                          | 1161      | 100.0   | 100.0         |                    |

**d2 ocupatia dvs. actuala (principala)**

|                                  | Frequency | Percent | Valid Percent | Cumulative Percent |
|----------------------------------|-----------|---------|---------------|--------------------|
| Valid 1 agricultor               | 259       | 22.3    | 22.3          | 22.3               |
| 2 muncitori (meserias)           | 247       | 21.3    | 21.3          | 43.6               |
| 3 tehnician, maistru, functionar | 74        | 6.4     | 6.4           | 50.0               |
| 4 ocupatii cu studii superioare  | 106       | 9.1     | 9.1           | 59.1               |
| 6 elev, student                  | 52        | 4.5     | 4.5           | 63.6               |
| 7 pensionar                      | 267       | 23.0    | 23.0          | 86.6               |
| 8 casnica                        | 62        | 5.3     | 5.3           | 91.9               |
| 9 acum sunt somer                | 82        | 7.1     | 7.1           | 99.0               |
| 10 patron                        | 12        | 1.0     | 1.0           | 100.0              |
| Total                            | 1161      | 100.0   | 100.0         |                    |

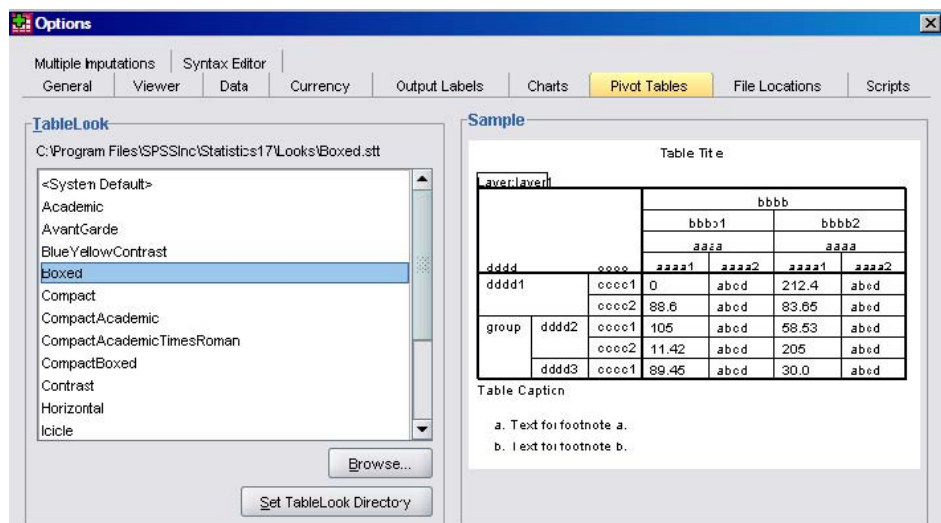
În tabul **Pivot Tables** (figura 3.3a) prefer să selectez, în secțiunea **TableLook**, opțiunea **Boxed**. Tabelul va avea toate celulele delimitate prin borduri. Acest lucru ne ajută, mai ales la tabelele mari, adică cu multe rânduri și coloane, să citim mai bine informația conținută de acestea. În figura 3.3b se observă diferența față de modul de prezentare a tabelului din figura 3.2b. Această opțiune este utilă în timpul analizelor. În rapoarte, articole, cărți sau alte materiale, nu vom copia tabelele din SPSS ca atare, ci le vom realiza în programul de editare a textului pe care îl folosim. Multe tabele oferite de SPSS conțin informații ce nu trebuie prezentate ca atare pentru cititor, acestea fiind utile în special analistului. De aceea, aceste informații trebuie eliminate sau prezentate în altă formă în cadrul materialului. Putem învăța să realizăm tabele ușor de citit, dacă parcurgem câteva articole publicate în jurnalele academice din domeniul care ne preocupă. O regulă de bază este : un tabel simplu este ușor de citit. Dacă acesta conține însă informații mai tehnice, atunci punem o notă imediat sub tabel în care explicăm cititorului cum trebuie să citească.

În tabul **File Locations**, la secțiunea **Startup Folders for Open and Save Dialogs**, prefer să bifez **Last folder used** (figura 3.4). O bază de date în format SPSS sau un fișier creat în acest program pot fi deschise fie dând dublu click pe fișier, fie din meniul **File > Open > Data**. O analiză poate dura mai multe zile, în funcție de complexitatea sa. În a doua zi de lucru, optez pentru a doua variantă de deschidere a fișierului. Dacă bifez **Last folder used**, atunci, mergând în **File > Open > Data**, programul ne va duce la ultimul fișier utilizat în ultima sesiune de lucru în acest program. Acest lucru este util pentru cei care au multe fișiere pe computer și, printre acestea, unul dedicat analizelor statistice, fișierul respectiv fiind astfel mult mai rapid de găsit la nevoie. Tot în acest tab, în secțiunea **Session Journal**, ne asigurăm că sunt bifate opțiunile **Record syntax in Journal** și **Append**.

Programul înregistrează toate operațiunile noastre. Dacă pierdem sintaxa, ștergând-o din greșeală, o vom găsi în jurnalul pe care îl ține SPSS. Acest jurnal poate fi salvat în fișierul predefinit de program sau într-un altul, ales de noi.

**Figura 3.3.** Tabul Pivot Tables (Edit > Options): modificarea designului tabelelor

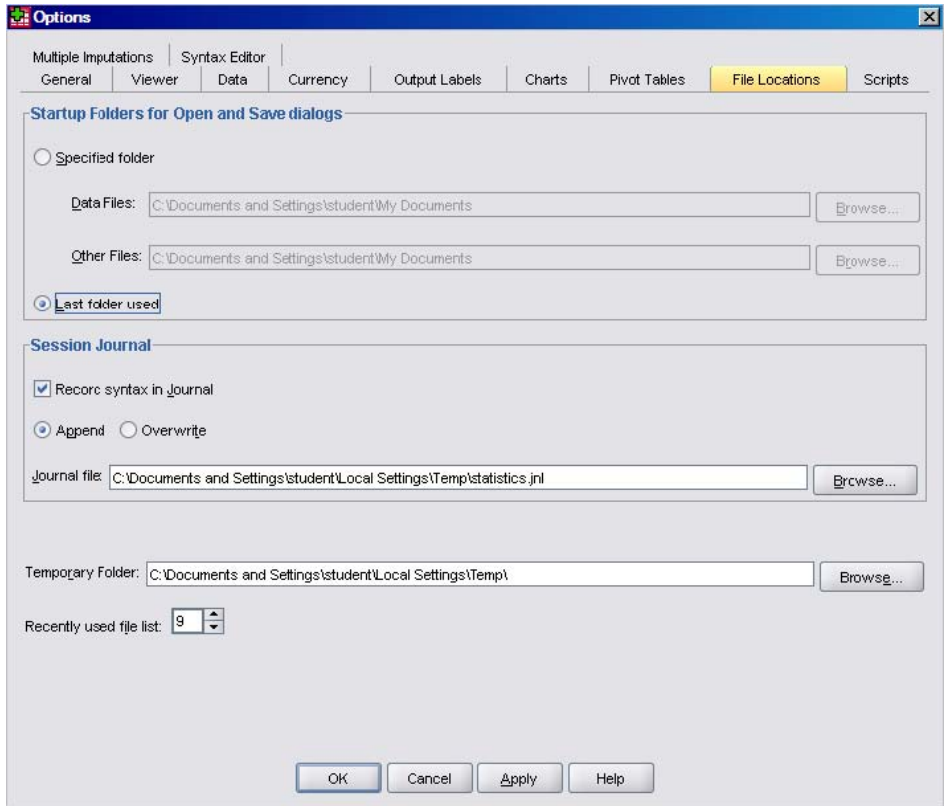
(a)



(b)

**d2 ocupatia dvs. actuala (principala)**

|       |                                  | Frequency | Percent | Valid Percent | Cumulative Percent |
|-------|----------------------------------|-----------|---------|---------------|--------------------|
| Valid | 1 agricultor                     | 259       | 22.3    | 22.3          | 22.3               |
|       | 2 muncitori (meserias)           | 247       | 21.3    | 21.3          | 43.6               |
|       | 3 tehnician, maistru, functionar | 74        | 6.4     | 6.4           | 50.0               |
|       | 4 ocupatii cu studii superioare  | 106       | 9.1     | 9.1           | 59.1               |
|       | 6 elev, student                  | 52        | 4.5     | 4.5           | 63.6               |
|       | 7 pensionar                      | 267       | 23.0    | 23.0          | 86.6               |
|       | 8 casnica                        | 62        | 5.3     | 5.3           | 91.9               |
|       | 9 acum sunt somer                | 82        | 7.1     | 7.1           | 99.0               |
|       | 10 patron                        | 12        | 1.0     | 1.0           | 100.0              |
|       | Total                            | 1161      | 100.0   | 100.0         |                    |

**Figura 3.4.** Tabul File Locations (Edit > Options) : fișierul de lucru și jurnalul SPSS

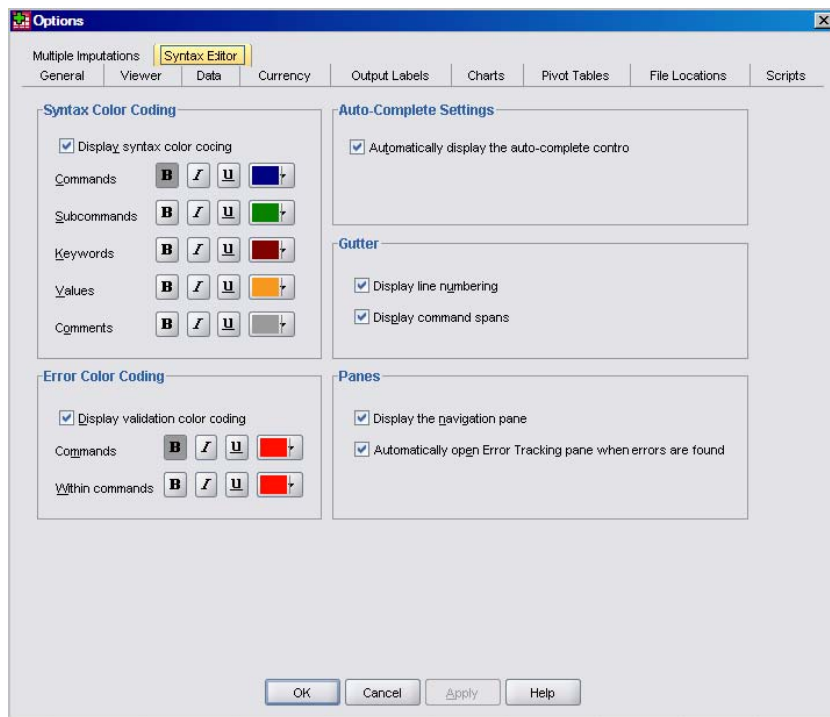
În fine, în tabul **Syntax Editor** (figura 3.5a) putem modifica modul în care arată sintaxa atunci când lucrăm cu ea. În ceea ce mă privește, îmi plac culorile alese implicit de program. Mi se pare extrem de utilă opțiunea de completare automată a unei comenzi care poate fi activată bifând **Automatically display the auto-complete control** (figura 3.5b). Dacă dorim să realizăm un tabel de frecvență știind că începutul comenzii este „fre...”, tastăm „fre...” și ni se va deschide fereastra din care putem alege comanda corectă. Această opțiune este foarte utilă pentru învățarea comenzilor uzuale. culorile pe care le afișează editorul sintaxei. Eu modific doar culoarea comentariilor, preferând un gri mai închis. În rest, sunt mulțumit de opțiunile implicite ale programului.

Fereastra sintaxei are două secțiuni : în partea din stânga se află lista comenzilor, iar în partea din dreapta sunt toate comenzile care, rulate, ne vor da analizele dorite. Lista comenzilor ne ajută să navigăm prin sintaxă, când aceasta conține multe comenzi. Atunci când, din greșeală, am scris o sintaxă greșită, dacă am bifat opțiunea **Automatically open Error Tracking pane when errors are found**, atunci, după cum se observă în figura 3.5c, se deschide o a treia secțiune care ne indică rândul unde se găsește eroarea, comanda care conține eroarea și informații

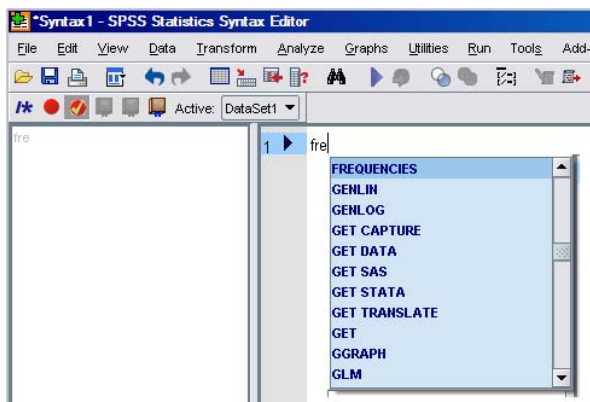
despre eroare. În acest caz, am introdus greșit numele variabilei pentru care dorim să realizăm tabelul de frecvență : variabila v1 nu există în baza de date.

**Figura 3.5.** Tabul Syntax Editor (Edit > Options) : cum putem face sintaxa mai ușor de utilizat

(a)

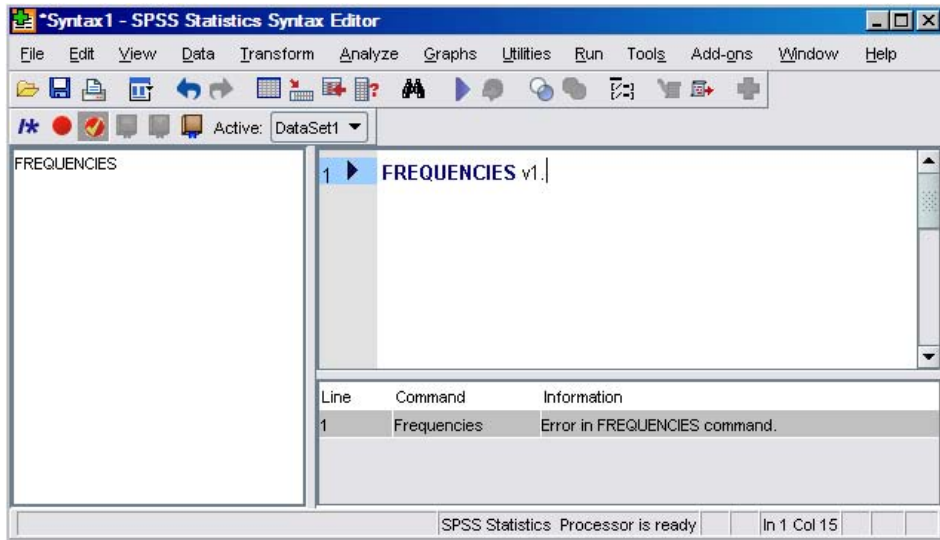


(b)





(c)



Am prezentat și opțiunile de la sintaxă pentru utilizatorii începători entuziaști. Pentru a parcurge această lucrare nu este nevoie să lucrăm cu sintaxa. Pe măsură ce vă familiarizați cu programul ar fi util să începeți să utilizați și sintaxa pe lângă meniuri. Veți constata că timpul acordat unei analize se reduce considerabil. În plus, veți avea jurnalul întregii analize la care puteți reveni oricând pentru referințe. O parte dintre sintaxele aferente comenzilor utilizate în această lucrare pot fi consultate pe pagina de internet dedicată acestora.

## 3.2. Pe scurt, despre structura programului :

### Data și Variable View

Aceste două elemente sunt esențiale în program : **Data View**, respectiv **Variable View**. Probabil că ați înțeles care este diferența dintre ele.


#### 3.2.1. Data View

**Data View** este secțiunea unde putem vizualiza datele. Dacă baza de date conține informații culese prin aplicarea unui chestionar, atunci fiecare rând va reprezenta un chestionar, iar fiecare coloană va reprezenta o variabilă. O celulă conține înregistrarea informației pentru un singur individ cu privire la o variabilă. Dacă variabila conține informații despre vârstă, atunci celula pentru rândul 1 va conține vârsta individului de pe rândul 1. În figura 3.6a este prezentată o secțiune din

baza de date DCV 2010. Rândul 1 reprezintă un român cărui i-au fost puse întrebările din chestionarul DCV 2010. Cifra 1 din dreptul variabilei nrchest pentru acest rând reprezintă numărul alocat acestui chestionar de către cercetător. Observăm că în figura 3.6a este selectată coloana nrchest. Cifra 1 din dreptul rândului 1 și coloanei d1 reprezintă sexul respondentului la chestionarul cu numărul 1. De unde știm ce reprezintă d1? Dar cifra 1? Vom afla proprietățile acestei variabile în **Variable View**. Să mai zăbovim puțin asupra interfeței **Data View**.


**Figura 3.6.** Data View

(a)



|   | nrchest | d1 | d2 | d3   | d4 |
|---|---------|----|----|------|----|
| 1 | 1       | 1  | 2  | 2006 | 1  |
| 2 | 2       | 0  | 1  | 99   | 2  |
| 3 | 3       | 0  | 1  | 1955 | 2  |
| 4 | 4       | 0  | 7  | 97   | 97 |
| 5 | 5       | 1  | 1  | 99   | 2  |

(b)



|   | nrchest | d1       | d2              | d3   | d4             | d5      |
|---|---------|----------|-----------------|------|----------------|---------|
| 1 | 1       | masculin | muncitori (...) | 2006 | salariat       | industr |
| 2 | 2       | feminin  | agricultor      | 99   | pe cont pro... | agric   |
| 3 | 3       | feminin  | agricultor      | 1955 | pe cont pro... | agric   |
| 4 | 4       | feminin  | pensionar       | 97   | 97             |         |
| 5 | 5       | masculin | agricultor      | 99   | pe cont pro... | agric   |

O parte dintre meniurile SPSS vă sunt familiare pentru că includ comenzi pe care le folosiți în mod frecvent în alte aplicații software uzuale. De exemplu, meniul **File** ne permite să deschidem documente, dar și să le salvăm. Meniul **Edit** ne permite să copiem (*copy*) și să lipim (*paste*) diferite elemente. Meniul **Window** ne permite să aranjăm documentele deschise astfel încât să le vizualizăm cât mai pe placul nostru. Meniul **Help** conține o mulțime de informații care ne ajută să înțelegem mai bine programul. Acest meniu se păstrează și când trecem în **Variable View**. În tabelul 3.1 sunt prezentate unele dintre cele mai utilizate comenzi, care vor fi discutate în această carte.

**Tabelul 3.1.** Meniuri frecvent utilizate

| Meniu       | Submeniu                                  | Utilitate  |
|-------------|---|--|
| <b>File</b> | New Data<br>New Syntax<br>New Output      | Realizăm o bază de date, un fișier de sintaxă sau unul de output, fără informații în ele.  |
|             | Open Data<br>Open Syntax<br>Open Output   | Deschidem o bază de date, un fișier de sintaxă sau unul de output care conțin informații.  |
|             | Save<br>Save as                           | Salvăm fișierele pe măsură ce lucrăm.<br>Salvăm fișierele sub alt nume sau în alt loc pe computer.   |
|             | Recently Used Data<br>Recently Used Files | Putem deschide un fișier cu care am lucrat într-o sesiune anterioară, fără a-l mai căuta pe computer.  |
| <b>Edit</b> | Insert Variable                           | Putem introduce manual o variabilă căreia îi definim, ulterior, proprietățile.   |
|             | Go To Case                                | Putem să găsim rapid un rând din baza de date.   |
|             | Go To Variable                            | Putem să găsim rapid o variabilă din baza de date, dacă îi știm numele.  |
|             | Options                                   | Putem să setăm programul conform preferințelor personale.  |
| <b>View</b> | Status Bar                                | Activăm sau dezactivăm <b>Status Bar</b> .   |
|             | Value Labels                              | Putem să vizualizăm în <b>Data View</b> etichetele atribuite codurilor (figura 3.6b).  |
|             | Variables                                 | Trecem din ecranul <b>Data View</b> în ecranul <b>Variable View</b> .  |
| <b>Data</b> | Identify Duplicate Cases                  | Putem verifica dacă, după unul sau mai multe criterii, am introdus în baza de date de mai multe ori același caz. Acest lucru se poate întâmpla, de exemplu, când chestionarele sunt aplicate prin e-mail și același respondent ne trimite chestionarul său de pe două adrese de e-mail diferite. |
|             | Sort Cases                                | Ordonăm cazurile în ordine crescătoare sau descrescătoare în funcție de una sau mai multe variabile. Putem să ordonăm și variante combinate.   |
|             | Merge Files                               | Unim două baze de date. Putem uni două baze care conțin aceleași cazuri, însă cu variabile diferite, dar și două baze care conțin cazuri diferite, însă cu aceleași variabile.   |
|             | Split File                                | Separăm baza de date după un criteriu. Analiza rulată este prezentată în același output comparativ pe grupurile definite de criteriul respectiv.   |
|             | Select Cases                              | Activăm sau dezactivăm anumite cazuri astfel încât să rulăm analizele doar pe anumite unități. Putem crea baze de date, pornind de la cea inițială.  |
|             | Weight Cases                              | Ponderăm baza de date. În prealabil, trebuie realizată variabila de ponderare.   |

|                   |  |   |
|-------------------|--|---|
| <b>Transform</b>  | Compute  | Realizăm o variabilă nouă, folosind o formulă și/ sau o funcție predefinită de SPSS.  |
|                   | Recode Into Same Variable                        | Modificăm codurile unei variabile, dar fără a-i modifica structura inițială.  |
|                   | Recode Into Different Variables                  | Modificăm structura unei variabile din baza de date. Rezultatul este o variabilă nouă.  |
| <b>Analyze</b>    | Descriptive statistics > Frequencies             | Realizăm tabele de frecvență, calculăm diferiți indicatori ai tendinței centrale, ai variației și/sau ai poziționării și creăm grafice.   |
|                   | Descriptive statistics > Explore                 | Explorăm datele. Putem testa asumptia distribuției normale folosind indicatori statistici și grafice.   |
|                   | Descriptive statistics > Crosstabs               | Realizăm tabele de contingență, inclusiv testul de semnificație <i>chi square</i> (hi-pătrat). Calculăm diferiți indicatori de asociere între variabile categoriale. Putem crea și graficul specific încrucișării variabilelor categoriale.               |
|                   | Descriptive statistics > P-P Plots sau Q-Q Plots | Testăm grafic abaterea de la distribuția normală.   |
|                   | Compare means > One-Sample T Test                | Comparăm media unei variabile din baza de date cu media furnizată de cercetător.  |
|                   | Compare means > Independent-Sample T Test        | Comparăm mediile a două grupuri.  |
|                   | Compare means > One-Way ANOVA                    | Comparăm mediile a cel puțin trei grupuri.  |
|                   | Correlate > Bivariate                            | Corelăm două variabile metrice.   |
|                   | Correlate > Partial                              | Corelăm două variabile metrice, controlând altă variabilă.  |
| <b>Regression</b> | Regression > Linear                              | Explicăm variația unei variabile metrice (dependentă), folosind simultan mai mulți predictori : rulăm analiza de regresie liniară.  |
|                   | Regression > Curve Estimation                    | Verificăm dacă între două variabile metrice există o relație liniară.   |
| <b>Graphs</b>     |  | Realizăm grafice.   |
| <b>Window</b>     | Split  | Putem împărți imaginea în <b>Data View</b> , astfel încât să vizualizăm datele în cel puțin două secțiuni. În figura 3.7 este prezentată împărțirea implicită activată prin utilizarea meniului. A nu se confunda cu meniul <b>Data &gt; Split File</b> . |
| <b>Help</b>       |  | Permite accesul la informații detaliate despre capacitățile programului.  |

Am enumerat în tabelul 3.1 informațiile pe care un începător trebuie să le acumuleze rapid. După ce acesta le-a înțeles, iar utilizarea lor este deja o rutină, tranziția către analizele mai complicate devine mult mai ușoară.

**Figura 3.7.** Meniul Window > Split : rezultatul împărțirii

|   | nrchest | d1 | nrchest | d1 | d2 |
|---|---------|----|---------|----|----|
| 1 | 1       |    | 1       | 1  | 2  |
| 2 | 2       |    | 2       | 0  | 1  |
| 3 | 3       |    | 3       | 0  | 1  |
| 4 | 4       |    | 4       | 0  | 7  |
| 5 | 5       |    | 5       | 1  | 1  |
| 6 | 6       |    | 6       | 0  | 1  |
| 7 | 7       |    | 7       | 1  | 1  |
| 2 | 2       |    | 2       | 0  | 1  |
| 3 | 3       |    | 3       | 0  | 1  |
| 4 | 4       |    | 4       | 0  | 7  |
| 5 | 5       |    | 5       | 1  | 1  |
| 6 | 6       |    | 6       | 0  | 1  |

### 3.2.2. Variable View

În acest meniu creăm variabilele și le definim proprietățile. Spre deosebire de **Data View**, în acest caz, rândul este o variabilă, iar coloanele sunt proprietăți diferite ale acestuia.

**Figura 3.8.** Variable View

|   | Name    | Type    | Width | Decimals | Label               | Values           | Missing | Columns | Align | Measure |
|---|---------|---------|-------|----------|---------------------|------------------|---------|---------|-------|---------|
| 1 | nrchest | Numeric | 8     | 0        | numarul chesi...    | None             | None    | 8       | Right | Scale   |
| 2 | d1      | Numeric | 8     | 0        | sexul               | {0, feminin}...  | None    | 8       | Right | Nominal |
| 3 | d2      | Numeric | 8     | 0        | ocupatia dvs. a...  | {1, agriculto... | None    | 8       | Right | Nominal |
| 4 | d3      | Numeric | 8     | 0        | din ce an aveti ... | None             | None    | 8       | Right | Scale   |
| 5 | d4      | Numeric | 8     | 0        | statutul ocupati... | {1, salariat}... | None    | 8       | Right | Nominal |
| 6 | d5      | Numeric | 8     | 0        | domeniul de sc...   | {1, agricultu... | None    | 8       | Right | Nominal |

Orice variabilă are un nume (coloana **Name**), căruia îi atribuim o etichetă (coloana **Label**). În figura 3.8 observăm, de exemplu, că variabila nrchest are eticheta „numarul chesi...”, iar variabila d1 are eticheta „sexul”, pe când variabila d2 are eticheta „ocupatia dvs. a...”. Putem vedea eticheta întreagă, adică „numarul chestionarului” sau „ocupatia dvs. actuala (principala)” dacă mergem cu mouse-ul între **Label** și **Values** și tragem de linia care le separă. Observăm că eticheta nu folosește diacritice. Uneori, în funcție și de setările computerului pe care este deschisă baza de date, acestea nu sunt citite corect, fiind înlocuite cu un simbol cum ar fi semnul de întrebare. De aceea, prefer să am două variante ale bazei de date : una cu diacritice și una fără diacritice.

Variabilele care au coduri, cum sunt aici sexul (d1), ocupația (d2) sau statutul ocupațional (d4), trebuie etichetate. Acest lucru se face în coloana **Values**. Atunci

când codurile nu sunt etichetate în dreptul variabilei respective, în coloana **Values** apare textul **None**.

În SPSS introducem, de regulă, numere. De aceea, majoritatea variabilelor vor fi numerice (coloana **Type**). Dacă introducem text, atunci tipul se schimbă în **String**.

Pentru că variabilele vizibile în figura 3.8 nu au valori cu zecimale, atunci în coloana **Decimals** ne asigurăm că avem valoarea 0. Dacă o variabilă are valori cu o zecimală, vom înlocui 0 cu 1, iar dacă are valori cu două zecimale, vom înlocui 0 cu 2 ș.a.m.d.

Ar fi indicat ca fiecărei variabile să îi fie definit corect nivelul de măsurare în coloana **Measure**. Astfel, vom beneficia de ajutor suplimentar din partea SPSS care, în anumite meniuri, dacă nivelul de măsurare este definit corect, va sugera diferite modalități de lucru.

Aș mai menționa aici doar coloana **Missing** în care instruim programul, introducând codurile aferente, pentru ignorarea nonrăspunsurilor în analize.

### 3.3. Ponderarea bazei de date (Data > Weight Cases)


În acest volum pornesc de la asumția că datele disponibile sunt culese prin utilizarea unui design de eșantionare probabilist. Un eșantion este probabilist atunci când toate obiectele care fac parte din populația de referință a studiului au o șansă diferită de zero de a fi selectate în eșantion (Levy și Lemeshow, 2008). Folosesc cuvântul „obiect”, pentru că, în funcție de nevoile de cercetare, putem fi interesați să extragem un eșantion de persoane (români adulți, cu vârsta egală sau mai mare de 18 ani sau elevi din clasele I-VIII, care fac parte din școli în care a fost implementat un program de reducere a abandonului școlar sau sunt consumatori ai iaurtului cu fructe produs de o anumită companie etc.), dar și de lucruri (mașini produse de o anumită companie care ies de pe linia de producție într-o lună, ouă care provin din găini crescute la sol și ouă care provin din găini crescute în baterii etc.). În toate exemplele fac referire la eșantioane de persoane.

Pentru a extrage un eșantion probabilist, avem nevoie de un cadru de eșantionare. Să presupunem că vrem să extragem un eșantion de persoane adulte cu vârsta de 18 ani și peste, neinstituționalizate. Designul frecvent utilizat în România este cel de tip stratificat, multistadial, cu selecție aleatoare în fiecare stadiu. După ce sunt selectate localitățile, se aleg secțiile de vot și, în final, cei care vor fi intervievați din cadrul fiecărei secții alese anterior. Informațiile despre distribuția populației României, grupată în funcție de regiunile de dezvoltare, ariile culturale (Sandu, 1999) sau regiunile istorice încrucișate cu mărimea orașelor și tipul de sat (aparținător sau reședință de comună), mărimea satelor după numărul de locuitori sau gradul de dezvoltare al localităților rurale pot fi culese de la Institutul Național de Statistică<sup>1</sup>. Informații

---

1. Institutul Național de Statistică : <http://www.insse.ro>.

despre secțiunile de vot și membrii acestora pot fi culese de la primăriile localităților selectate sau de la Autoritatea Electorală Permanentă<sup>1</sup>. Dacă aceste informații sunt actualizate și armonizate corespunzător, iar operatorii de teren respectă instrucțiunile primite de la cercetători, atunci structura eșantionului pentru variabile-cheie cum ar fi mediul de rezidență, sexul, vârsta, educația ș.a. ar fi similară cu structura populației. În practică, există situații când cele două structuri nu se suprapun perfect. De aceea, se recurge la ponderare. Programul de statistică este instruit să ia în considerare într-o măsură mai mare ceea ce este subreprezentat în eșantion și într-o măsură mai mică ceea ce este suprareprezentat în eșantion. Acest lucru se face prin construirea unei variabile denumită pondere (*weight*). De exemplu, în cercetarea World Values Survey din 2012 (WVS 2012), al cărui chestionar a fost aplicat și în România, a fost folosită o variabilă de ponderare. Calcularea ponderilor este un proces destul de laborios care nu face obiectul acestei cărți. O descriere detaliată, într-un context comparativ, poate fi consultată pe platforma ESS EduNet<sup>2</sup> pusă la dispoziție în cadrul proiectului European Social Survey<sup>3</sup> sau în lucrările dedicate eșantionării, cum ar fi cea scrisă de Levy și Lemeshow (2008) care, în *Sampling of Populations: Methods and Applications*, dedică un întreg capitol construirii ponderilor, acesta fiind scris de Paul S. Biemer și Sharon L. Christ. De asemenea, vă recomand să consultați materialele metodologice ale unor cercetări cum ar fi European Values Study<sup>4</sup> sau European Quality of Life Survey<sup>5</sup>.

Echipa din România care a aplicat a creat o variabilă de ponderare pentru baza de date. Înainte de a începe analizele propriu-zise, baza de date trebuie ponderată. Există situații, ca aceasta la care mă refer, în care baza de date pe care trebuie să o folosim pentru a rula anumite analize statistice nu este creată de noi. Analistul primește baza de date pregătită pentru analiză. Aceasta ar trebui să conțină și variabila de ponderare. În WVS 2012, căutând în **Variable View**, am aflat că variabila de ponderare este V258. Căutarea în **Variable View** se poate face în mai multe moduri. Putem să navigăm, derulând vertical în coloana **Name** și coloana **Label**, căutând cuvintele-cheie „pondere” sau „weight” sau un alt cuvânt asemănător. Mai rapid ar fi să dăm click în prima celulă din coloana **Label**. Apoi apăsăm iconița , cu care suntem familiarizați din alte programe utilizate în viața de zi cu zi. Sau putem apăsa simultan tastele CTRL + F. Se deschide fereastra din figura 3.9. În secțiunea **Find** tastăm cuvântul „weight”. Am ales acest cuvânt-cheie pentru că baza de date WVS 2012 este etichetată în limba engleză, fiind o cercetare comparativă la nivel internațional, iar datele din

1. Autoritatea Electorală Permanentă : <http://www.roaep.ro>.

2. ESS EduNet : <http://essedunet.nsd.uib.no/cms/topics/weight>.

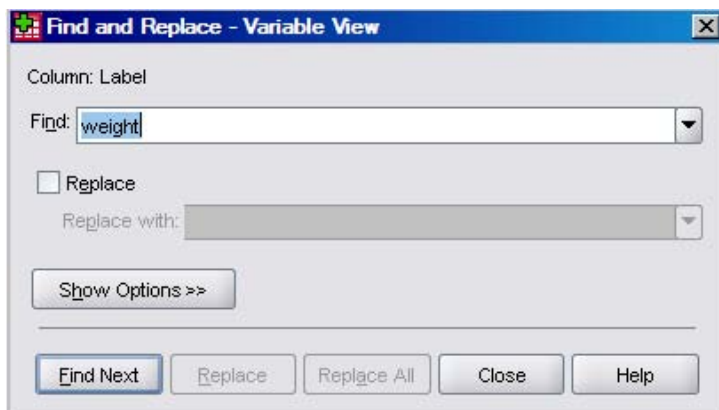
3. European Social Survey : <http://www.europeansocialsurvey.org>.

4. European Values Study : <http://www.europeanvaluesstudy.eu>.

5. European Quality of Life Survey : <http://www.eurofound.europa.eu/surveys/eqls/index.htm>.

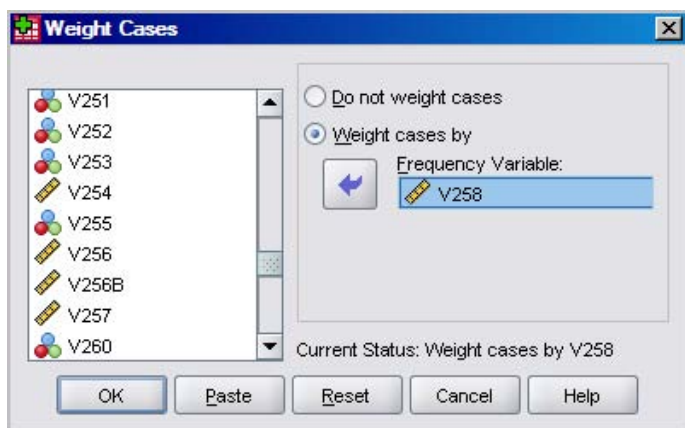
România sunt integrate în același fișier cu datele din alte țări. Apăsăm butonul **Find Next** o dată sau de mai multe ori, până când găsim ceea ce căutăm. Dacă etichetele ar fi fost scrise în limba română, am fi folosit cuvântul-cheie „pondere”. Nu există o regulă: cuvintele sunt alese în funcție de ce vrem să găsim. Când nu suntem siguri cu privire la forma sub care este folosit cuvântul, tastăm doar o parte din acesta: „weig” sau „pond”.

**Figura 3.9.** Find : căutare după un cuvânt-cheie



Ponderarea se face din meniul **Data > Weight Cases**. În figura 3.10 este prezentată fereastra cu modificările efectuate, pregătită doar pentru a apăsa butonul **OK**. Căutăm variabila V258 în lista de variabile din stânga. Inițial, este selectată opțiunea **Do not weight cases**. Dacă baza de date nu trebuie ponderată, această opțiune rămâne neschimbată. Aici bifăm opțiunea **Weight cases by**. S-a activat secțiunea **Frequency Variable** în care introducem, folosind săgeata, variabila V258.

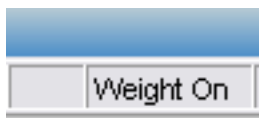
**Figura 3.10.** Data > Weight Cases : meniul în care activăm ponderea





În output nu se va întâmpla nimic, adică nu va fi produs nici un tabel sau un grafic. În baza de date, indiferent că ne aflăm în meniul **Data View** sau **Variable View**, iar opțiunea **Status Bar** este activată, ar trebui să fie, în colțul din dreapta jos, afișată expresia **Weight On**, ca în figura 3.11.

**Figura 3.11.** Confirmare vizuală că ponderea este activă



Ponderea poate fi dezactivată în același meniu, **Data > Weight Cases**, selectând **Do not weight cases** și apăsând **OK**. Revenind în meniul **Data View** sau **Variable View**, ar trebui să fi dispărut confirmarea **Weight On** prezentată în figura 3.11.

Observăm în figura 3.10 că variabila de ponderare pe care o solicită SPSS este **Frequency Variable**. Acest lucru înseamnă că valorile pe care le ia variabila de ponderare sunt numere de tipul 1, 2, 3, 100, 130 etc. În **Help**, de altfel și când rulăm diferite analize, suntem avertizați că, atunci când cazurile primesc pondere egală cu zero sau ponderi cu numere negative (cu minus), acestea sunt eliminate din analiză. Unele analize acceptă și ponderi de tipul 1.2, 0.7 etc., iar alte analize nu acceptă deloc ponderi. Trebuie să vă documentați bine înainte de a rula o analiză pe o bază ponderată, pentru a vedea în ce măsură este corectată structura eșantionului.

### 3.4. Filtrarea bazei de date (**Data > Select Cases**)

A filtra o bază de date înseamnă a selecta din total doar cazurile care îndeplinesc unul sau mai multe criterii. Cazurile care îndeplinesc criteriul de filtrare rămân active în baza de date, iar celelalte sunt dezactivate. De asemenea, putem să le copiem într-o bază de date diferită. Mai putem să ștergem din baza de date inițială cazurile care nu satisfac criteriul respectiv.

Cercetarea World Values Survey presupune aplicarea unui chestionar cu multe întrebări comune în mai multe țări într-o perioadă dată de timp. Cercetătorii doresc să compare țările respective după caracteristicile măsurate în chestionar. După încheierea muncii de teren, vor exista atâtea baze de date câte țări au fost incluse în cercetare. Aceste baze de date sunt unite într-un singur fișier. De exemplu, în cazul acestei cercetări, puteți descărca baza de date care conține toate țările și toate etapele din perioada 1981-2005 de pe site-ul World Values Survey, iar, în curând, acesteia îi va fi adăugată și ultima etapă care, în România, s-a derulat în 2012. Așadar, avem o bază de date care conține atât eșantionul românesc, cât și pe cel german, dar și altele. Să presupunem că suntem interesați să lucrăm

doar cu eșantionul românesc. Pentru că multe cazuri și multe variabile înseamnă o bază de date mare în termeni de dimensiuni (megabiți), acest lucru s-ar putea traduce prin durate mai mari de procesare a analizei solicitate computerului pe care lucrați. Dacă acesta ne permite să lucrăm cu volume mari de date, ne-ar putea totuși interesa și partea estetică – să vizualizăm doar variabilele și cazurile care ne interesează. În oricare dintre aceste contexte, vom utiliza un filtru în baza de date integrată care ne permite să extragem o nouă bază, care să conțină doar eșantionul românesc. Dacă nu vrem să avem mai multe baze de date pe computer, va trebui doar să activăm un filtru care va instrui programul să ia în considerare doar cazurile ce ne interesează, iar după încheierea activității care solicita filtrul, îl vom dezactiva și vom vizualiza, din nou, baza inițială cu toate cazurile.

### 3.4.1. Activarea unui filtru : lucrăm pe baza de date inițială

Să presupunem că vrem doar să păstrăm active anumite cazuri, fără a crea o bază distinctă. Lucrăm doar cu datele culese în România.

Vrem să rulăm o analiză doar pentru bărbați : Care este procentul bărbaților români care se declară fericiți sau foarte fericiți ?

Mai întâi, trebuie să găsim variabila care indică sexul respondenților. Căutând în **Variable View**, aflăm că aceasta se numește V240. Pentru a activa un filtru în baza de date, trebuie să cunoaștem valorile (codurile) variabilei/variabilelor care constituie filtrul respectiv. Care este codul bărbaților? Pentru a răspunde la această întrebare, realizăm un tabel de frecvență (tabelul 3.2) folosind meniul **Analyze > Descriptive Statistics > Frequency**. Pentru a vedea codurile, trebuie să fi făcut modificările în meniul **Edit > Options** așa cum le-am discutat ceva mai devreme. Codul bărbaților (*Male*) este 1.

**Tabelul 3.2.** Tabel de frecvență : Care sunt codurile folosite pentru bărbați și pentru femei ?

| V240 Sex |          |           |         |               |                    |
|----------|----------|-----------|---------|---------------|--------------------|
|          |          | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid    | 1 Male   | 723       | 48.1    | 48.1          | 48.1               |
|          | 2 Female | 780       | 51.9    | 51.9          | 100.0              |
|          | Total    | 1503      | 100.0   | 100.0         |                    |

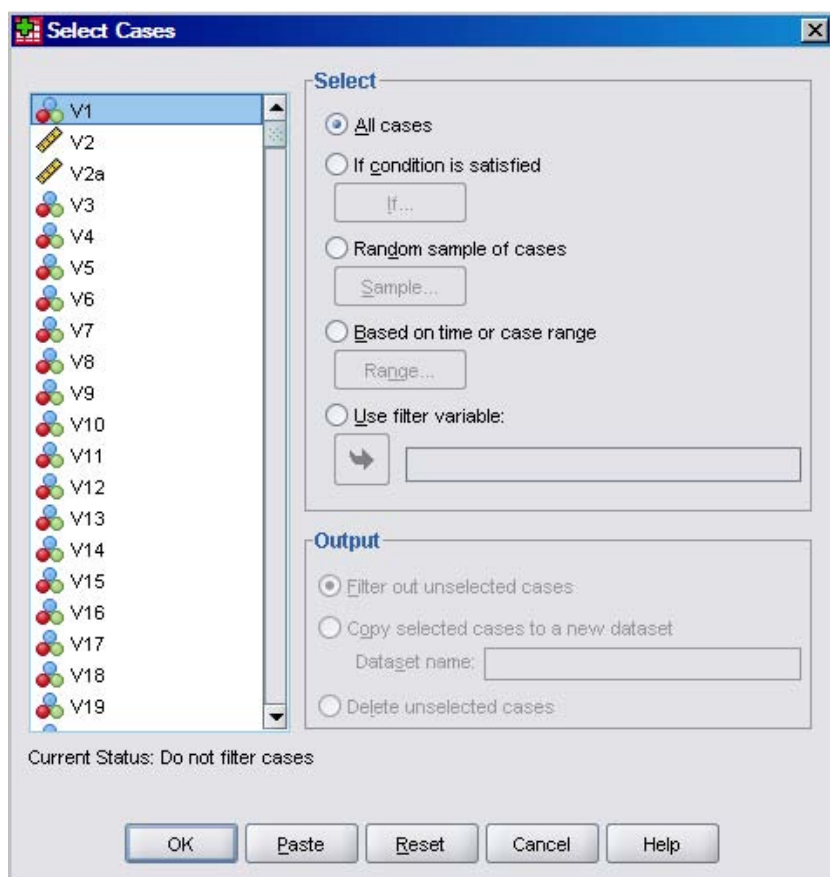
Pentru a selecta doar bărbații, trebuie să utilizăm filtrul : **V240 = 1**. Astfel, vor rămâne activi în baza de date doar bărbații. Să reținem structura filtrului : **numele variabilei = cod**. Pentru a activa acest filtru, mergem în meniul **Data > Select Cases** (Figura 3.12a). Inițial, în secțiunea **Select**, este bifat **All cases**. SPSS utilizează, în această situație, toate cazurile din baza de date. Pentru a activa filtrul dorit, trebuie să bifăm **If condition is satisfied**. Observăm că se activează

butonul **If**. După ce am apăsât butonul **If**, se deschide fereastra în care vom pune condiția prin care instruiem SPSS să păstreze activi doar bărbații (figura 3.12b). Căutăm variabila V240 în lista de variabile din stânga și, folosind săgeata, o trecem în secțiunea din dreapta sus. Apoi introducem filtrul : **V240 = 1**. Bărbații trebuie să rămână activi. Apăsăm **Continue**. Ne asigurăm că în fereastra inițială (figura 3.12a), în secțiunea **Output**, este bifată opțiunea **Filter out unselected cases**. Apăsăm **OK**.

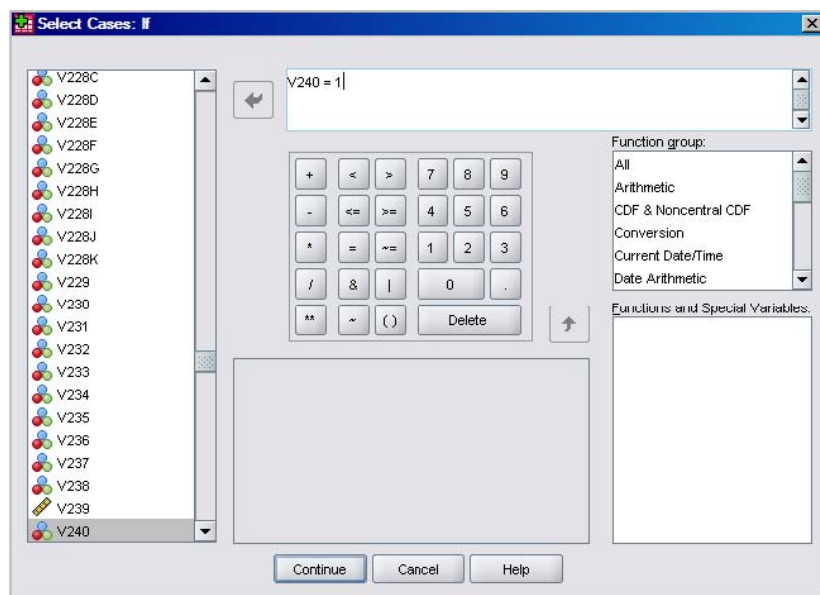
Începătorii cred că ar trebui să utilizeze butoanele pe care le oferă programul (figura 3.12b). Apăsând, de exemplu, semnul „=” și cifra 1, programul ne va pune automat și spațiile necesare între elementele distincte. Semnele „<=” și „>=” înseamnă „mai mic sau egal”, respectiv „mai mare sau egal”. Semnul „~=” înseamnă „diferit de”. Semnul „&” înseamnă „și”, iar semnul „|” înseamnă „sau”.

**Figura 3.12.** Meniul Data > Select Cases : fereastra inițială prin care activăm, dezactivăm, copiem sau ștergem cazuri

(a)



(b)



Filtrul este activ. Trebuie să verificăm corectitudinea operațiunii efectuate. În **Data View** vedem că unele rânduri nu sunt tăiate (2, 3, 6, 7 și 10), iar altele sunt tăiate (1, 4, 5, 8 și 9) (tabelul 3.3a). Acesta este modul programului SPSS de a ne spune că un filtru este activat. Dar, pentru a fi siguri că filtrul activ este corect, în această situație, realizăm un tabel de frecvență pentru variabila folosită în filtru, V240 (tabelul 3.3b). Observăm că doar bărbații sunt activi, deci filtrul activ este cel dorit.

**Tabelul 3.3.** Tabel de frecvență : verificarea corectitudinii filtrului

(a)

|              | V1 | V2  |
|--------------|----|-----|
| <del>1</del> | 6  | 242 |
| 2            | 6  | 242 |
| 3            | 6  | 242 |
| <del>4</del> | 6  | 242 |
| <del>5</del> | 6  | 242 |
| 6            | 6  | 242 |
| 7            | 6  | 242 |
| <del>8</del> | 6  | 242 |
| <del>9</del> | 6  | 242 |
| 10           | 6  | 242 |


(b)

| V240 Sex |        |           |         |               |                    |
|----------|--------|-----------|---------|---------------|--------------------|
|          |        | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid    | 1 Male | 723       | 100.0   | 100.0         | 100.0              |

Nu ne rămâne decât să realizăm un alt tabel de frecvență pentru variabila care ne arată procentul bărbaților români fericiți sau foarte fericiți. Această variabilă poartă numele V10. Folosind meniul **Analyze > Descriptive Statistics > Frequencies**, obținem tabelul 3.4, unde observăm că 13 % sunt „foarte fericiți” (*Very happy*) și 58 % sunt „destul de fericiți” (*Rather happy*)<sup>1</sup>. Citim procentele valide (**Valid Percent**) care sunt calculate din totalul bărbaților care și-au declarat nivelul de fericire, adică au răspuns la V10. Acest total este 719 bărbați, spre deosebire de totalul general care este 723 de bărbați. Folosind procentele cumulate (**Cumulative Percent**), puteam să spunem că 71 % dintre bărbații români se declarau foarte fericiți sau destul de fericiți în 2012.

**Tabelul 3.4.** Tabel de frecvență : Distribuția fericirii  
în rândul bărbaților români (WVS 2012)

| V10 Feeling of happiness |                    |           |         |               |                    |
|--------------------------|--------------------|-----------|---------|---------------|--------------------|
|                          |                    | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                    | 1 Very happy       | 91        | 12.6    | 12.7          | 12.7               |
|                          | 2 Rather happy     | 418       | 57.7    | 58.1          | 70.7               |
|                          | 3 Not very happy   | 184       | 25.5    | 25.6          | 96.4               |
|                          | 4 Not at all happy | 26        | 3.6     | 3.6           | 100.0              |
|                          | Total              | 719       | 99.4    | 100.0         |                    |
| Missing                  | -2 No answer       | 1         | .2      |               |                    |
|                          | -1 Don't know      | 3         | .4      |               |                    |
|                          | Total              | 4         | .6      |               |                    |
|                          | Total              | 723       | 100.0   |               |                    |

În unele situații, filtrele de care avem nevoie sunt mai complexe. Folosind aceleași date, dorim să aflăm care este nivelul de fericire al bărbaților care au educație superioară. Filtrul include acum două variabile : sexul și educația. Mai întâi, trebuie să aflăm care sunt variabilele de care avem nevoie pentru analiză. Știm că sexul este V240. Educația este V248. Fericirea este V10. Să ne amintim : am găsit numele variabilelor în **Variable View**, dând click într-o celulă în coloana **Label**, apăsând pe iconița reprezentând binoclu  și tastând „sex”, „educ” sau „happ”. Pasul următor presupune să aflăm codurile pe care le vom folosi pentru a crea filtrul. Pentru realizarea acestui obiectiv trebuie să alcătuim un tabel de frecvență pentru fiecare dintre cele două variabile de filtrare, sexul (V240) și educația (V248). Deja știm codurile pentru

1. Traducerea în limba română este preluată din chestionarul românesc al WVS 2012.

sex, așa că, folosind meniul **Analyze > Descriptive Statistics > Frequencies**, realizăm unul doar pentru educație (tabelul 3.5). Dacă, anterior, am creat alte tabele de frecvență și nu am închis baza de date, veți remarca faptul că în fereastra meniului există acele variabile. Pentru a reveni la setările inițiale din meniu, apăsăm butonul

Reset

. Codurile pentru educație superioară sunt 8 și 9. Dacă nu am fi avut etichete pentru coduri, nu am fi știut care dintre acestea reprezintă educația superioară. Aici este folosită o schemă de clasificare a nivelurilor educaționale care permite comparația între țări. Puteți consulta, de exemplu, *International Standard Classification of Education (ISCED)*<sup>1</sup> pentru a înțelege mai bine această idee.

**Tabelul 3.5.** Tabel de frecvență : Care sunt codurile pentru bărbați și pentru mediul rural ?

| <b>V248 Highest educational level attained</b> |  |           |         |               |                    |
|--|--|-----------|---------|---------------|--------------------|
|  |  | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid  | 1 No formal education                                      | 5         | .7      | .7            | .7                 |
|  | 2 Incomplete primary school                                | 13        | 1.8     | 1.8           | 2.5                |
|  | 3 Complete primary school                                  | 26        | 3.6     | 3.7           | 6.2                |
|  | 4 Incomplete secondary school: technical/ vocational type  | 27        | 3.7     | 3.7           | 9.9                |
|  | 5 Complete secondary school: technical/ vocational type    | 188       | 26.0    | 26.3          | 36.2               |
|  | 6 Incomplete secondary school: university-preparatory type | 120       | 16.6    | 16.8          | 52.9               |
|  | 7 Complete secondary school: university-preparatory type   | 165       | 22.8    | 23.1          | 76.0               |
|  | 8 Some university-level education, without degree          | 52        | 7.2     | 7.3           | 83.3               |
|  | 9 University-level education, with degree                  | 119       | 16.5    | 16.7          | 100.0              |
|  | Total  | 715       | 98.8    | 100.0         |                    |
| Missing  | –2 No answer   | 8         | 1.2     |               |                    |
|  | Total  | 723       | 100.0   |               |                    |

Așadar, filtrul poate fi scris sub forma : **V240 = 1 & (V248 = 8 | V248 = 9)**. Prima secțiune a filtrului, dinaintea semnului **&**, o cunoaștem : sunt bărbații. A doua secțiune se referă la educația superioară : observăm că, dacă folosim mai multe condiții pentru aceeași variabilă, trebuie să îi introducem numele de fiecare

1. <http://www.uis.unesco.org/Education/Pages/international-standard-classification-of-education.aspx>.

dată. Pentru că un respondent nu poate alege în chestionar la V248 decât un singur răspuns, trebuie să folosim semnul | (sau).

Pentru verificare realizăm un tabel de frecvență pentru fiecare dintre cele două variabile de filtrare (tabelul 3.6). Observăm că au rămas active în baza de date doar codurile pentru bărbații cu studii superioare.

**Tabelul 3.6.** Tabele de frecvență : verificarea corectitudinii  
filtrului V240 = 1 & (V248 = 8 | V248 = 9)

| V240 Sex |        |           |         |               |                    |
|----------|--------|-----------|---------|---------------|--------------------|
|          |        | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid    | 1 Male | 171       | 100.0   | 100.0         | 100.0              |

| V248 Highest educational level attained |   |           |         |               |                    |
|---|---|-----------|---------|---------------|--------------------|
|   |   | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                                   | 8 Some univer-<br>sity-level education,<br>without degree | 52        | 30.5    | 30.5          | 30.5               |
|   | 9 University-level<br>education, with degree              | 119       | 69.5    | 69.5          | 100.0              |
|   | Total   | 171       | 100.0   | 100.0         |                    |

Un filtru poate fi scris, uneori, în mai multe forme. Încercați acest lucru în situația dată.

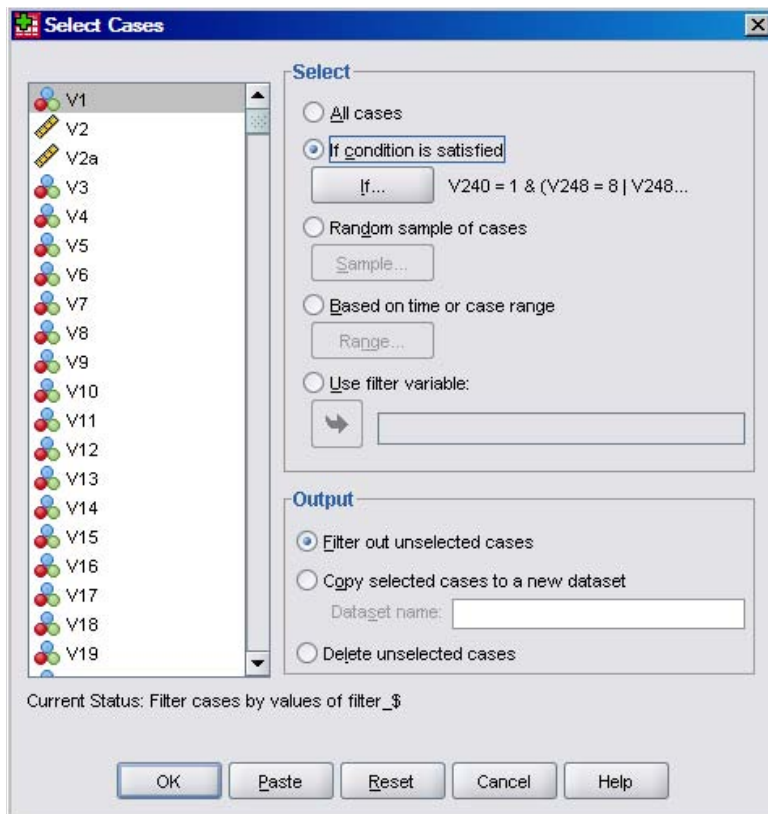
Acum putem rula analiza propriu-zisă : aflarea procentului bărbaților români cu studii superioare care se declară fericiți (tabelul 3.7). 86% dintre aceștia se declară foarte fericiți sau destul de fericiți. Procentele sunt calculate din totalul de răspunsuri valide, adică 170. Când filtrați baza de date, fiți atenți la cazurile care rămân active : dacă vă rămân puține cazuri, atunci trebuie să vă întrebați ce relevanță are analiza respectivă.

**Tabelul 3.7.** Tabel de frecvență : Distribuția fericirii în rândul bărbaților români cu studii superioare (WVS 2012)

| V10 Feeling of happiness |                    |           |         |               |                    |
|--------------------------|--------------------|-----------|---------|---------------|--------------------|
|                          |                    | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                    | 1 Very happy       | 19        | 10.8    | 10.9          | 10.9               |
|                          | 2 Rather happy     | 129       | 75.2    | 75.7          | 86.5               |
|                          | 3 Not very happy   | 22        | 12.9    | 13.0          | 99.5               |
|                          | 4 Not at all happy | 1         | .5      | .5            | 100.0              |
|                          | Total              | 170       | 99.4    | 100.0         |                    |
| Missing                  | -1 Don't know      | 1         | .6      |               |                    |
|                          | Total              | 171       | 100.0   |               |                    |

Filtrul rămâne activ până când îl dezactivăm. Dezactivarea se face din același meniu **Data > Select Cases**. Trebuie doar să bifăm **All cases** și apoi să apăsăm butonul **OK**. Atunci când filtrul este activ, în **Data View** sau **Variable View**, în colțul din dreapta jos observăm pe **Status Bar** **Filter On** **Weight On**. După ce am bifat **All cases** și am apăsut **OK**, va dispărea **Filter On** din **Status Bar**. În încheiere, să observăm fereastra meniului cu toate modificările efectuate (figura 3.13). Observăm în dreapta butonului **If** condiția activă și, sub lista de variabile și deasupra butoanelor, expresia **Current Status : Filter cases by values of filter\_\$**. Această expresie ne indică faptul că SPSS a creat o variabilă care ia valorile 1 și 0, unde 1 este codul atribuit cazurilor care îndeplinesc condiția și 0, codul celor care nu o îndeplinesc. Dacă dorim să reutilizăm filtrul fără a mai face toate aceste operațiuni, atunci putem redenumi această variabilă în **Variable View** și, apoi, când avem nevoie de ea, o putem introduce în secțiunea **Use filter variable**. Dacă nu o redenumim, data viitoare când creăm un filtru folosind condiții noi, aceasta va fi eliminată și vom pierde informația inițială.

**Figura 3.13.** Meniul **Data > Select Cases** : fereastra cu filtrul care menține active doar anumite cazuri

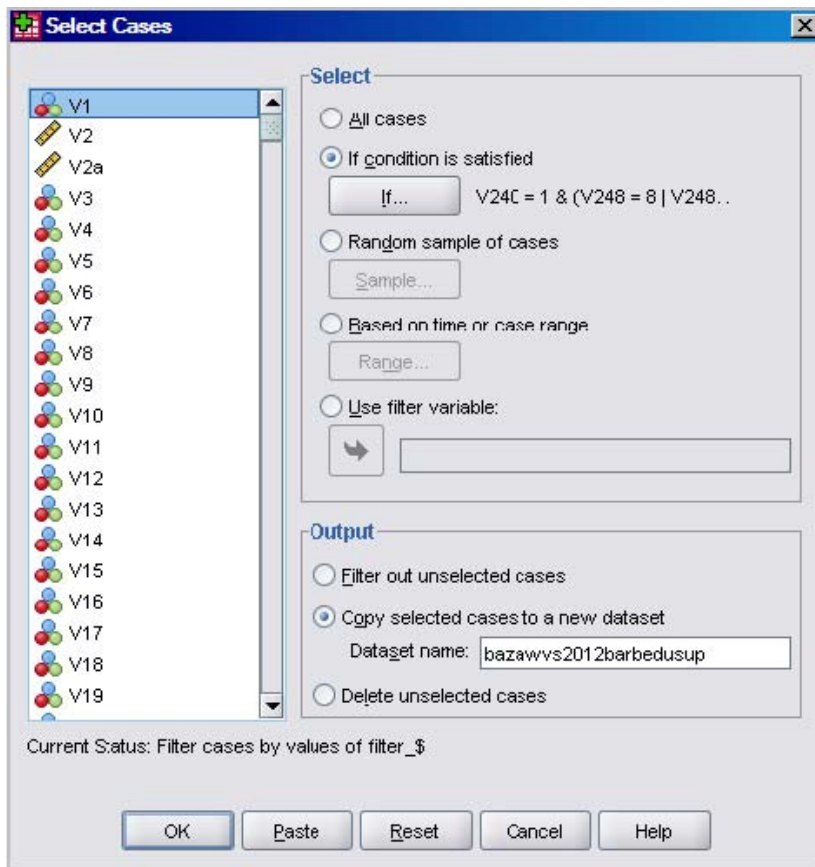




### 3.4.2. Crearea unei baze de date folosind un filtru

Folosind acest meniu, putem crea și o bază de date care conține doar cazurile ce îndeplinesc anumite condiții. Singurul lucru pe care trebuie să îl facem este ca, în loc să bifăm **Filter out unselected cases**, să bifăm **Copy selected cases to a new dataset** (figura 3.14).

**Figura 3.14.** Meniul Data > Select Cases : fereastra cu filtrul care creează o bază nouă de date



Când bifăm **Copy selected cases to a new dataset** se activează opțiunea **Dataset name**. Aici trebuie să introducem un nume pentru noua bază de date, care trebuie să respecte condițiile impuse numelor variabilelor : să înceapă cu o literă și să nu conțină spații între caractere. Ar fi de preferat să fie și scurt. Apăsând **OK**, SPSS creează o bază de date care trebuie salvată pe computer, această bază conținând doar cazurile pe care le definește filtrul.

În figura 3.14 se observă că mai avem, în secțiunea **Output**, opțiunea **Delete unselected cases**. Aceasta este utilă doar dacă ați salvat baza de date originală și lucrați pe o copie a acesteia. În caz contrar, veți pierde informații greu de recuperat după această acțiune distructivă.

### 3.5. Separarea bazei de date (Data > Split File)

Utilizarea filtrelor este un lucru obișnuit în manipularea bazei de date și în analiza datelor din aceasta.

În unele situații dorim să comparăm rezultatul unei analize pentru două sau mai multe grupuri. Care este procentul bărbaților foarte fericiți prin comparație cu cel al femeilor fericite? Predictorii fericirii aleși în cazul femeilor și în cel al bărbaților funcționează la fel? Ideea de bază este că, prin separarea bazei de date (**split file**), putem vizualiza outputul unei analize pentru două sau mai multe grupuri distincte. Pentru aceasta pot fi folosite și filtre, ceea ce este o chestiune de gust, în multe situații.

SPSS ne permite să separăm baza de date în funcție de o variabilă categorială care conține cel puțin două grupuri, cum ar fi bărbați *versus* femei, locuitori din mediul rural *versus* locuitori din mediul urban, români *versus* germani *versus* bulgari, căsătoriți *versus* divorțați *versus* văduvi etc. Variabila categorială este variabila de separare. Alte variabile vor fi utilizate pentru a rula o analiză pentru fiecare dintre aceste grupuri. În tabelul 3.8 este prezentat tabelul de frecvență al variabilei fericire pentru bărbați, respectiv, femei.

Lucrăm, așadar, cu două tipuri de variabile : cea de separare și cea sau cele pe care le folosim în analize statistice. Aici am separat în funcție de sex și am făcut o analiză statistică pentru fericire. Pentru situația de față, este mai util să realizăm un tabel de contingență, despre care vom vorbi în alt capitol al acestui volum. Mi se pare mai utilă această opțiune pe care ne-o oferă SPSS atunci când rulăm un model multivariat, cum ar fi o regresie liniară multiplă. Dacă presupunem că modelul funcționează diferit pentru bărbați și pentru femei, atunci putem vedea rezultatul în output în funcție de opțiunea de separare prezentată aici.

Revenind la principiul opțiunii, variabila de separare va fi, întotdeauna, o variabilă categorială : sex, mediu de rezidență, stare civilă etc. Putem folosi și variabile metrice, cum ar fi vârsta, doar dacă aceasta a fost recodificată în prealabil : 18-34, 35-64, 65+ . Atunci când recodificăm, creând categorii, trebuie ca numărul de cazuri ce revine fiecărei categorii să fie rezonabil de mare astfel încât să aibă sens comparația dintre grupurile rezultate. De aceea, vom folosi variabile de separare cu puține categorii, mai ales când eșantioanele sunt mici ca volum.

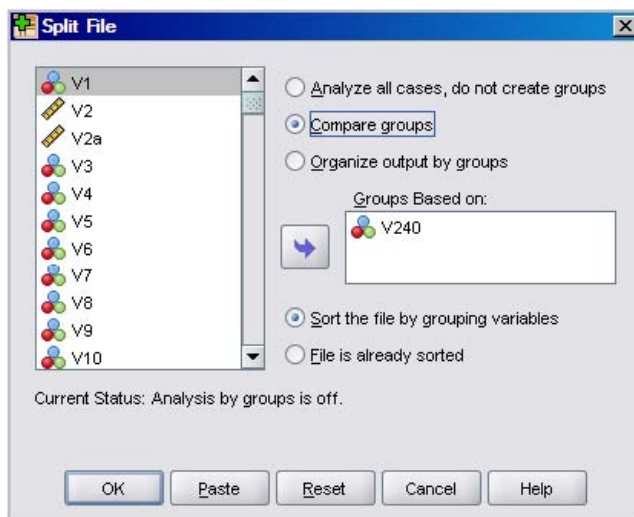
**Tabelul 3.8.** Reprezentare grafică a rezultatului separării bazei de date

| <b>V10 Feeling of happiness</b> |         |                    |           |         |               |                    |
|---------------------------------|---------|--------------------|-----------|---------|---------------|--------------------|
| V240 Sex                        |         |                    | Frequency | Percent | Valid Percent | Cumulative Percent |
| 1 Male                          | Valid   | 1 Very happy       | 91        | 12.6    | 12.7          | 12.7               |
|                                 |         | 2 Rather happy     | 418       | 57.7    | 58.1          | 70.7               |
|                                 |         | 3 Not very happy   | 184       | 25.5    | 25.6          | 96.4               |
|                                 |         | 4 Not at all happy | 26        | 3.6     | 3.6           | 100.0              |
|                                 |         | Total              | 719       | 99.4    | 100.0         |                    |
|                                 | Missing | –2 No answer       | 1         | .2      |               |                    |
|                                 |         | –1 Don't know      | 3         | .4      |               |                    |
|                                 |         | Total              | 4         | .6      |               |                    |
|                                 |         | Total              | 723       | 100.0   |               |                    |
|                                 |         |                    |           |         |               |                    |
| 2 Female                        | Valid   | 1 Very happy       | 114       | 14.6    | 14.7          | 14.7               |
|                                 |         | 2 Rather happy     | 415       | 53.3    | 53.5          | 68.2               |
|                                 |         | 3 Not very happy   | 213       | 27.3    | 27.5          | 95.7               |
|                                 |         | 4 Not at all happy | 33        | 4.3     | 4.3           | 100.0              |
|                                 |         | Total              | 776       | 99.5    | 100.0         |                    |
|                                 | Missing | –2 No answer       | 3         | .4      |               |                    |
|                                 |         | –1 Don't know      | 1         | .1      |               |                    |
|                                 |         | Total              | 4         | .5      |               |                    |
|                                 |         | Total              | 780       | 100.0   |               |                    |
|                                 |         |                    |           |         |               |                    |

Dacă dorim să calculăm vârsta medie a bărbaților și a femeilor din România și să avem această informație într-un singur tabel, atunci putem folosi separarea. Variabila de separare va fi sexul, iar variabila pentru care calculăm media va fi vârsta. Variabila sex are numele V240, iar variabila vârstă are numele V242. Mai întâi, alcătuim câte un tabel de frecvență pentru a verifica dacă există nonrăspunsuri și pentru a ne familiariza cu cele două variabile. Variabila sex are două coduri, 1, pentru bărbat și 2, pentru femeie. Variabila vârstă are foarte multe valori, cea minimă fiind 18 ani și cea maximă fiind 85 de ani. La variabila sex nu există nonrăspunsuri. La variabila vârstă există trei persoane care nu și-au declarat vârsta. Observăm că, în baza de date, acestea au fost deja definite ca nonrăspunsuri, pentru că există în tabelul de frecvență secțiunea **Missing**, sub **Total**. Celor trei persoane care nu și-au declarat vârsta le-au fost atribuite codul –2, „nu răspund”. Aceste operațiuni fiind deja realizate, putem trece la analiza propriu-zisă.

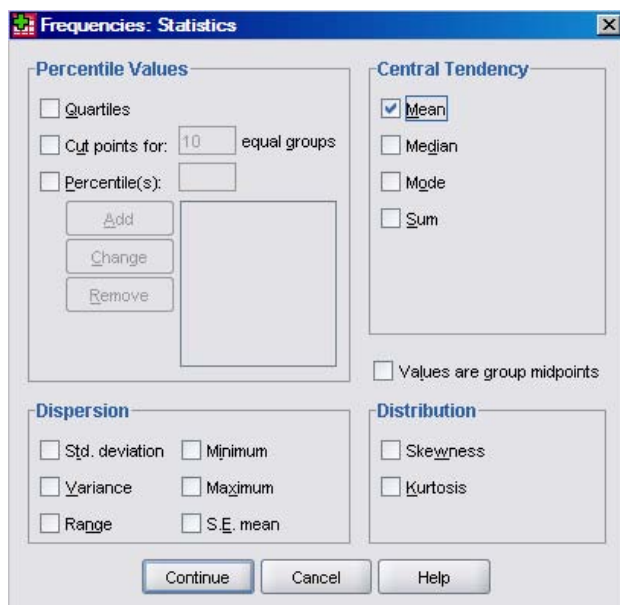
Mai întâi trebuie să separăm baza de date. Acest lucru se face mergând în meniul **Data > Split File** (figura 3.15). În fereastra care se deschide, inițial este bifată opțiunea **Analyze all cases, do not create groups**. Bifăm opțiunea **Compare groups**. Introducem variabila sex (V240) în secțiunea **Groups Based on**. Apăsăm **OK**. SPSS ne anunță că este activă opțiunea de separare a bazei de date în **Data View** sau **Variable View**, în colțul din dreapta jos pe **Staus Bar** : Weight On Split by V240 .

Figura 3.15. Meniul Data &gt; Split File



Acum putem calcula media vârstei. Acest lucru îl putem face din meniul cu care ne-am obișnuit deja, **Analyze > Descriptive Statistics > Frequencies**. De data aceasta, vom utiliza și butonul **Statistics** unde, în secțiunea **Central Tendency**, bifăm **Mean** (figura 3.16).

Figura 3.16. Meniul Analyze &gt; Descriptive Statistics &gt; Frequencies : cum calculăm media unei variabile



Rezultatul analizei este prezentat în tabelul 3.9. Cele două medii sunt 45 de ani pentru bărbați, respectiv 47 de ani pentru femei.

**Tabelul 3.9.** Media vârstei : tabel obținut prin separarea bazei de date

| Statistics |   |         |       |
|------------|---|---------|-------|
| V242 Age   |   |         |       |
| 1 Male     | N | Valid   | 721   |
|            |   | Missing | 2     |
|            |   | Mean    | 45.00 |
| 2 Female   | N | Valid   | 779   |
|            |   | Missing | 1     |
|            |   | Mean    | 47.40 |

Opțiunea de separare rămâne activă până când o dezactivați. Este o situație similară cu cea de la filtrare. Trebuie să intrați înapoi în meniul **Data > Split File** și să bifați opțiunea **Analyze all cases, do not create groups**. Când opțiunea de separare nu mai este activă, textul **Split by** din **Status Bar** dispăre.

### 3.6. Exerciții

Notă : exercițiile presupun utilizarea bazei de date European Values Study 2008 România, disponibilă gratuit la ZACAT – GESIS Online Study Catalogue<sup>1</sup>.

1. Este baza de date ponderată ? Dacă nu, ponderați baza de date.
2. Câți bărbați consideră că prietenii și cunoștințele lor sunt importanți în viață ? Aplicați un filtru, pentru a răspunde la întrebare.
3. Câte femei consideră că prietenii și cunoștințele lor sunt importanți în viață ? Aplicați un filtru, pentru a răspunde la întrebare.
4. Câte femei consideră că familia este importantă în viață ? Aplicați un filtru, pentru a răspunde la întrebare.
5. Câți bărbați consideră că familia este importantă în viață ? Aplicați un filtru, pentru a răspunde la întrebare.
6. Câți locuitori ai localităților cu peste 100.000 de locuitori consideră că religia este importantă ? Aplicați un filtru, pentru a răspunde la întrebare.
7. Câți locuitori ai localităților cu mai puțin de 100.000 de locuitori consideră că religia este importantă ? Aplicați un filtru, pentru a răspunde la întrebare.
8. Unde sunt mai mulți oameni fericiți : în localitățile cu mai puțin de 100.000 de locuitori sau în localitățile care au peste 100.000 de locuitori ? Separați (Split) baza de date, pentru a răspunde la această întrebare.

1. <http://zacat.gesis.org/webview/index.jsp>.

9. Cine discută mai frecvent despre politică : bărbații sau femeile ? Separați baza de date, pentru a răspunde la această întrebare.
10. În ce regiune de dezvoltare sunt cei mai mulți oameni fericiți ? Separați baza de date, pentru a răspunde la această întrebare.

## 4. Curățarea și validarea unei baze de date

Înainte de a trece la analiza datelor, trebuie să ne asigurăm că acestea nu conțin erori. Aici avem, de fapt, două idei. Una dintre ele este cea pe care o discutăm în acest capitol : eliminarea erorilor de culegere și de introducere a datelor. Acesta este procesul de curățare și de validare a bazei de date. A doua idee se referă la testarea validității și a fidelității măsurătorilor cu care lucrăm. În acest sens, putem consulta materiale cum ar fi cele scrise de Mărginean (1982), Saris și Gallhofer (2007) sau Carmines și Zeller (1979).

Curățarea și validarea unei baze de date constituie un pas esențial în procesul cercetării cantitative. Acesta este un proces pentru că toate activitățile specifice unei abordări cantitative a socialului sunt interconectate. Cel care primește sarcina să curețe baza de date va comunica permanent cu echipa care a coordonat activitatea de teren. Acesta poate să identifice erori în baza de date care trebuie verificate prin consultarea chestionarului. Curățarea nu este o activitate făcută într-un birou obscur de cineva care rulează coduri.

În zilele noastre, multe companii de cercetare nu mai tipăresc chestionarele pe hârtie, ci folosesc o metodă de înregistrare digitală. Tableta este un instrument foarte util în acest sens. Folosind această abordare, este redusă considerabil cantitatea de muncă și de resurse materiale, umane și temporale necesare pentru finalizarea cercetării.

Informațiile prezentate aici se aplică atât în situațiile în care realizați o cercetare proprie și parcurgeți toate etapele aferente, cât și în situațiile în care utilizați date culese și introduse într-o bază de date de altcineva. În a doua situație, teoretic, datele sunt deja curățate, iar baza este validată. În practică, însă, este bine să realizați propria verificare : în fond, scăpările altora afectează rezultatul analizelor dumneavoastră.

Etapele esențiale pentru curățarea și validarea bazei de date sunt :

- etichetarea variabilelor și valorilor variabilelor, acolo unde este necesar acest lucru ;
- dezactivarea nonrăspunsurilor ;
- verificarea introducerii eronate a unor coduri ;
- validarea logică prin urmărirea filtrelor din chestionar, dar și a unor întrebări factuale ;
- recodificarea unor variabile esențiale și construirea unor variabile noi.

Procesul de curățare ne ajută să apreciem posibilitatea de a utiliza sintaxa. Sintaxa este echivalentul în cod al clickurilor pe care le dați în meniuri. Din

sintaxă puteți rula chiar și comenzi care nu se regăsesc în meniuri. Sintaxa are mai multe avantaje, dintre care aș puncta : (a) avem un jurnal al operațiunilor pe care le-am realizat în bază, putând reveni oricând la ele pentru a le consulta sau a le rula pe o bază curată ; (b) scade timpul petrecut cu diferite operații. Nu trebuie să învățați comenzile. Pe unele dintre ele, pe măsură ce le utilizați, le veți reține fără probleme. Printre acestea se numără cele pentru tabelul de frecvență (**frequencies**), tabelul de contingență (**crosstabs**), recodificări (**recode**), pentru realizarea de noi variabile (**compute**) etc. Mai mult, SPSS ne oferă în toate meniurile butonul **Paste** care, apăsând după ce am bifat toate opțiunile dorite, le transformă în coduri pe care le putem salva și rula oricând.

## 4.1. Etichetarea variabilelor și a valorilor variabilelor

Am importat baza de date în SPSS și s-a deschis fereastra **Variable View** (figura 4.1). Trebuie să completăm informații pentru fiecare variabilă (fiecare rând) în coloanele **Label**, **Values** și **Missing**. SPSS ghidează analistul în anumite situații, sugerându-i analizele și graficele pe care le poate face în funcție de nivelul de măsurare a variabilelor selectate. Acest lucru este posibil dacă selectăm corect opțiunile din coloana **Measure**. Însă în practică, aceasta este o opțiune pe care o putem ignora, pentru că, pe măsură ce învățăm să lucrăm cu datele cantitative și avem mai multe cunoștințe de statistică, putem decide singuri în situațiile respective. Este chiar preferabil să controlați acțiunile programului, și nu să îl lăsați să ia decizii în locul dumneavoastră.

**Figura 4.1.** Variable View : baza de date înainte și după etichetare

|   | Name    | Type    | Width | Decimals | Label | Values | Missing | Columns | Align | Measure |
|---|---------|---------|-------|----------|-------|--------|---------|---------|-------|---------|
| 1 | nrchest | Numeric | 8     | 0        |       | None   | None    | 8       | Right | Scale   |
| 2 | d1      | Numeric | 8     | 0        |       | None   | None    | 8       | Right | Scale   |
| 3 | d2      | Numeric | 8     | 0        |       | None   | None    | 8       | Right | Scale   |
| 4 | d3      | Numeric | 8     | 0        |       | None   | None    | 8       | Right | Scale   |
| 5 | d4      | Numeric | 8     | 0        |       | None   | None    | 8       | Right | Scale   |
| 6 | d5      | Numeric | 8     | 0        |       | None   | None    | 8       | Right | Scale   |

|   | Name    | Type    | Width | Decimals | Label               | Values           | Missing | Columns | Align | Measure |
|---|---------|---------|-------|----------|---------------------|------------------|---------|---------|-------|---------|
| 1 | nrchest | Numeric | 8     | 0        | numarul chesti...   | None             | None    | 8       | Right | Scale   |
| 2 | d1      | Numeric | 8     | 0        | sexul               | {0, feminin}...  | None    | 8       | Right | Nominal |
| 3 | d2      | Numeric | 8     | 0        | ocupatia dvs. a...  | {1, agricultu... | None    | 8       | Right | Nominal |
| 4 | d3      | Numeric | 8     | 0        | din ce an avet ...  | None             | None    | 8       | Right | Scale   |
| 5 | d4      | Numeric | 8     | 0        | statutul ocupati... | {1, salariat}... | None    | 8       | Right | Nominal |
| 6 | d5      | Numeric | 8     | 0        | domeniul de ac...   | {1, agricultu... | None    | 8       | Right | Nominal |

În coloana **Label**, introducem explicații detaliate despre ce scrie în coloana **Name**. Dacă nu facem acest lucru, atunci când redeschidem baza de date nu vom ști ce înseamnă nrchest, d1, d2 etc. Nu putem găsi nume intuitive pentru toate variabilele. Chiar dacă avem mereu la îndemână un chestionar când lucrăm (de fapt, îl avem), nu este tocmai intuitiv să nu avem etichete în baza de date când



rulăm diferite analize. De regulă, în coloana **Label**, se trece chiar întrebarea din chestionar. Dacă este prea lungă, atunci o putem prescurta alegând cuvintele cele mai importante, astfel încât cei care vor lucra cu această bază de date să înțeleagă ușor informațiile respective.

Putem proceda în două moduri : (1) în meniul **Variable View** scriem în coloana **Label** în dreptul variabilei care ne interesează sau (2) deschidem un fișier de sintaxă și scriem comenzile, apoi le rulăm. Recomand varianta a doua pentru că ați putea șterge din greșală baza și nu veți mai avea sintaxa, caz în care va trebui să o luați de la capăt cu etichetarea și cu celelalte modificări din acest meniu. Un fișier de sintaxă poate fi creat din meniul **File > New > Syntax**. Se va deschide o fereastră nouă similară cu cea din figura 3.5b. Comanda prin care adăugăm o etichetă unui nume de variabilă, adică introducerea unei informații în coloana **Label**, este foarte simplă : **VARIABLE LABELS** sau, prescurtat, **VAR LAB**. Iată sintaxa pentru etichetarea celor șase variabile din figura 4.1 :

**VARIABLE LABELS** nrchest „numarul chestionarului”

**VARIABLE LABELS** d1 „sexul”

**VARIABLE LABELS** d2 „ocupatia dvs. actuala (principala)”

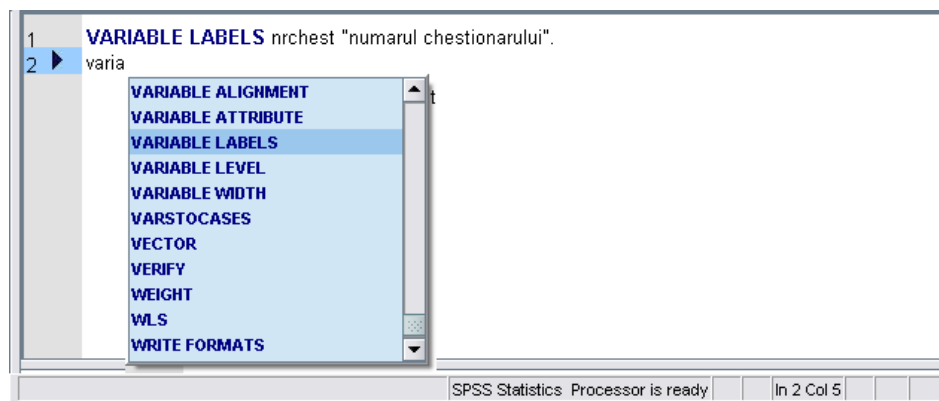
**VARIABLE LABELS** d3 „din ce an aveti aceasta ocupatie ? ”

**VARIABLE LABELS** d4 „statutul ocupational”

**VARIABLE LABELS** d5 „domeniul de activitate”

Pentru începători, găsesc utilă folosirea denumirii complete a comenzilor. Veți afla foarte rapid că puteți prescurta aceste comenzi. Acesta poate fi un exercițiu de familiarizare cu programul : care este varianta prin care puteți folosi doar o singură dată comanda **VARIABLE LABELS** pentru toate cele șase variabile ? Folosiți opțiunea **Help** a programului pentru a afla acest lucru.

**Figura 4.2.** Fișierul de sintaxă : afișarea listei derulante de comenzi



În figura 4.2 observăm că este suficient să tastăm primele litere din comandă și programul ne ajută deschizând o listă derulantă din care putem alege ceea ce ne interesează. Nimic mai simplu ! Puteți întreba : dar de unde știu care sunt

comenzile pe care trebuie să le folosesc? Lăsând acest volum la o parte, puteți căuta pe internet – SPSS are foarte mulți utilizatori – și printre aceștia se găsesc mulți entuziaști ai sintaxei. Cu siguranță, veți găsi rapid ceea ce doriți. Programul are un manual bogat care poate fi găsit în meniul **Help** sau chiar pe internet, pe pagina producătorului, IBM.

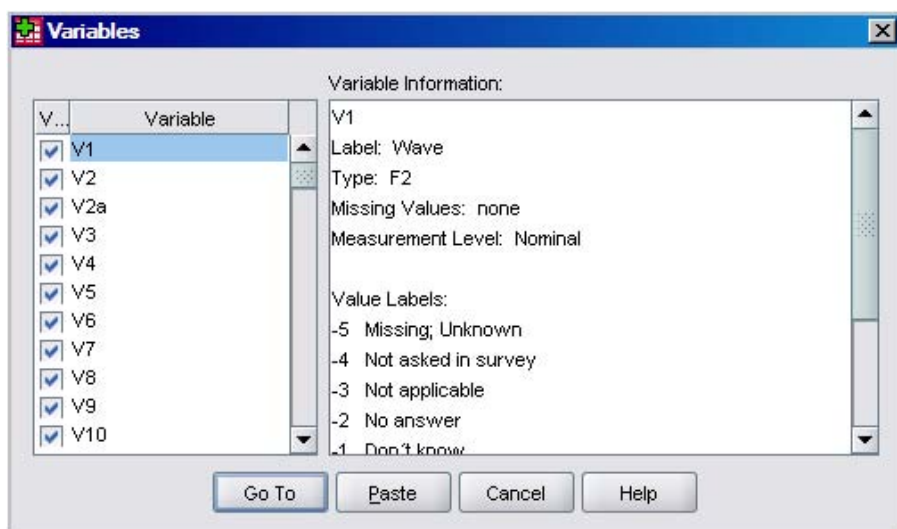
După ce am scris prima linie de sintaxă, în această situație, este suficient să selectăm rândul respectiv, **Copy** și **Paste** pe rândul următor. Modificăm nrchest cu d1 și, între ghilimele, scriem eticheta corespunzătoare. Apoi **Paste** pe rândul următor și modificăm nrchest cu d2 și, între ghilimele, scriem eticheta corespunzătoare. Repetăm până când am epuizat variabilele care trebuie etichetate.

Observăm următoarele :

- putem scrie comanda **VARIABLE LABELS** sau **VAR LAB** fie cu litere mici, fie cu MAJUSCULE. SPSS folosește în lista derulantă majuscule, dar acestea nu sunt obligatorii. Pentru a crește vizibilitatea în interiorul sintaxei, prefer să folosesc pentru comenzi majuscule, iar pentru comentarii litere mici.
- pe rând, între fiecare element al comenzii, lăsăm un spațiu, apăsând tasta spațiu. Comandă [spațiu] numele variabilei [spațiu] [ghilimele stânga] [eticheta] [ghilimele dreapta] [punct].
- după comanda **VARIABLE LABELS**, notăm numele variabilei, aici nrchest sau d1 sau d2 etc. SPSS oferă posibilitatea de a pune automat numele variabilei în fișierul de sintaxă. Putem merge în meniul **Utilities > Variables** (figura 4.3). Variabila nrchest este prima. Nu ne va fi de mare folos. Dar să presupunem că vrem să găsim rapid variabila V240 : dăm click în lista de variabile din stânga ferestrei pe oricare variabilă, astfel încât aceasta să fie selectată (de exemplu, aici este selectată V1). Apoi tastăm rapid primele două-trei caractere din numele variabilei care ne interesează, aici V240. Programul ne va duce imediat la variabila V240. Ne asigurăm că este selectată și apăsăm butonul **Go To**, dacă vrem să fie afișată în baza de date, sau butonul **Paste**, dacă vrem să fie copiată în sintaxă. Vom alege a doua opțiune. Acest meniu este foarte util atunci când variabilele nu au denumiri atât de intuitive ca d1, d2, V240 etc., ci mai greu de ținut minte, cum ar fi tvtot, trstlgl, prtvctbe etc., acestea fiind denumiri folosite în baza de date a cercetării *European Social Survey 2012*<sup>1</sup>.
- după numele variabilei, între ghilimele, scriem eticheta. Ghilimelele, în principiu, sunt necesare dacă eticheta conține caractere speciale cum ar fi cratima, semnul exclamării, punct etc. În plus, delimitează vizual sintaxa.
- întreaga comanda se încheie cu punct.

1. [http://www.europeansocialsurvey.org/docs/round6/survey/ESS6\\_appendix\\_a8\\_e01\\_0.pdf](http://www.europeansocialsurvey.org/docs/round6/survey/ESS6_appendix_a8_e01_0.pdf).

**Figura 4.3.** Meniul Utilities > Variables : cum găsim rapid o variabilă și cum îi copiem numele în sintaxă

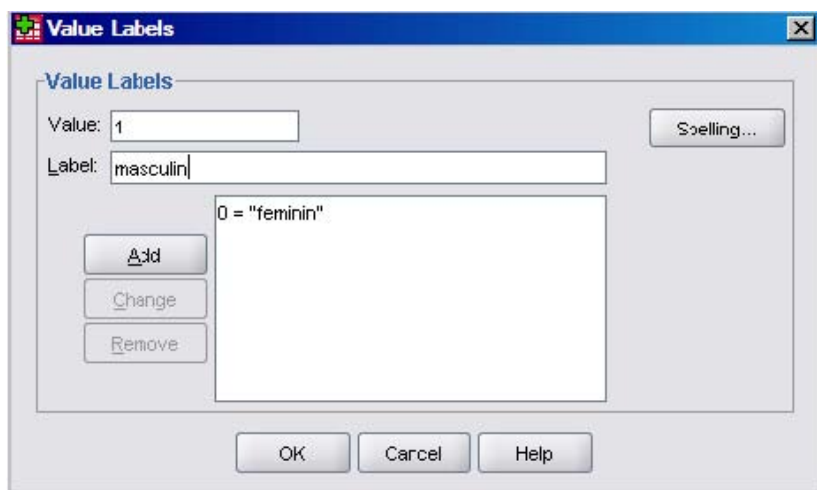


După ce am încheiat de notat sintaxa și am salvat fișierul, putem rula sintaxa. Deschidem meniul **Run**, unde există mai multe posibilități. Dacă vrem să rulăm doar o anumită comandă, și nu întregul fișier de sintaxă, atunci alegem **Selection**. Pentru a vedea modificările, mergem în **Variable View** (figura 4.1). Puteți rula sintaxa și fără să utilizați acest meniu : găsiți prescurtarea !

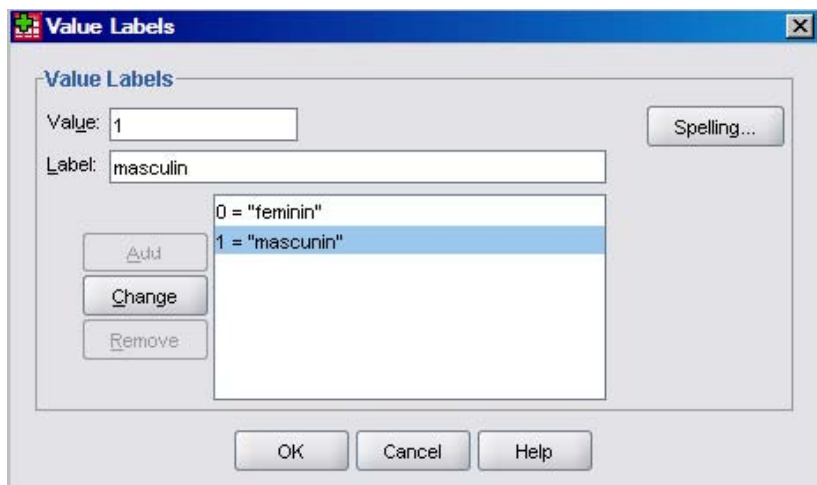
După ce am încheiat etichetarea variabilelor, trecem la etichetarea valorilor, acolo unde este cazul. Ceea ce discutăm acum se va finaliza prin introducerea unor informații în coloana **Values** din **Variable View**. Am precizat că nu tot timpul este nevoie să etichetăm valorile. Unele variabile au variante de răspuns care nu mai necesită explicații. Vârsta este măsurată în ani împliniți : știm ce înseamnă 46. Salariul din luna trecută este măsurat în lei : știm ce înseamnă 1.350. Numărul de camere pe care îl are locuința este... un număr : știm ce înseamnă 2. În schimb, alte variabile nu sunt măsurate la fel de intuitiv. Variantele de răspuns sunt exprimate numeric prin atribuirea unor coduri. Fericirea este măsurată prin întrebarea : „Luând în considerare toate aspectele vieții dvs., ați spune că sunteți... 1. Foarte fericit, 2. Destul de fericit, 3. Nu prea fericit, 4. Deloc fericit”. Respondentul alege o etichetă, dar în baza de date introducem codul. Dacă nu etichetăm codul, nu știm ce reprezintă acesta. Pentru precizarea ocupației, respondentul trebuie să aleagă dintre mai multe variante de răspuns : fiecare are un cod. Fiecare cod trebuie etichetat. Putem eticheta codurile fie în **Variable View**, fie în sintaxă, folosind o comandă simplă. Să ne ocupăm de prima variantă. Dăm click pe celula din dreptul variabilei dorite, aici d1, și al coloanei **Values**. Se vor activa trei puncte pe care dăm click. Se va deschide fereastra din figura 4.4a.

**Figura 4.4.** Variable View : etichetarea valorilor variabilei

(a)



(b)



Toate secțiunile sunt inițial goale. În celula **Value** introducem codul 0, iar în celula **Label** vom introduce eticheta „feminin”. Se va activa butonul **Add**, pe care îl apăsăm. Continuăm cu codul 1: în celula **Value**, introducem codul 1, iar în celula **Label** vom introduce eticheta „masculin”. Am putea avea și un cod de nonrăspuns. Pentru a verifica acest lucru trebuie să realizăm un tabel de frecvență pentru variabila d1 folosind meniul **Analyze > Descriptive statistics > Frequencies**. Dacă ar exista un cod de nonrăspuns, atunci ar trebui să îi alocăm și acestuia o etichetă. Dacă am tastat greșit, de pildă, „masculin”, atunci vom

selecta în celula mare eticheta scrisă greșit, vom modifica în celula **Label** și vom apăsa butonul **Change** (figura 4.4b).

Este mai rapid să utilizăm sintaxa, care este la fel de simplă ca cea utilizată la etichetarea variabilelor : **VALUE LABELS** sau **VAL LAB**. Pentru etichetarea variabilelor folosite ca exemplu aici, sintaxa va fi :

**VALUE LABELS d1**

0 „feminin”

1 „masculin”

**VALUE LABELS d2**

1 „agricultor”

2 „muncitor (meserias)”

3 „tehnician, maistru, functionar”

4 „ocupatie cu studii superioare”

5 „alta ocupatie”

6 „elev, student”

7 „pensionar”

8 „casnica”

9 „acum sunt somer”

10 „patron”

**VALUE LABEL d4**

1 „salariat”

2 „pe cont propriu”

3 „patron”

4 „zilier”

**VALUE LABEL d5**

1 „agricultura”

2 „industrie, constructii”

3 „transporturi, telecomunicatii”

4 „comerț, turism, intermediieri etc.”

5 „invatamant, cultura, cercetare, proiectare”

6 „sanatate”

7 „altele”

Structura sintaxei este aceeași ca la **VARIABLE LABELS**, cu diferența că etichetele și codurile sunt trecute pe rânduri separate.

În chestionar, la d5, varianta de răspuns cu codul 4 are o etichetă ceva mai lungă : „comerț, turism, intermediieri (financiare, imobiliare, pariuri etc.)”. SPSS permite un număr limitat de caractere pentru etichetele valorilor, de aceea am preferat să folosesc „etc.” în locul informației dintre paranteze. Dacă nu aș fi trunchiat eticheta, ar fi făcut-o SPSS, numai că într-un mod mai puțin intuitiv de citit. Aflați care este numărul maxim de caractere pe care le permite SPSS pentru etichetele valorilor.

În chestionarul DCV 2010 există mai multe variabile care au aceleași variante de răspuns, deci aceleași etichete. Putem folosi o singură comandă de etichetare a valorilor acestor variabile. Să luăm, de exemplu, variabilele d14-d27. Întrebarea din chestionar este : „În viața fiecăruia intervin o mulțime de condiții și împrejurări. Ele pot fi mai bune sau mai puțin bune. Mai jos sunt menționate o serie de asemenea aspecte. Vă rugăm să le caracterizați, în ceea ce vă privește, încercuind cifra corespunzătoare. Alegeți un singur răspuns la fiecare întrebare”. Variabilele cărora li se aplică această întrebare sunt :

|     |  | Foarte proastă(e) | Proastă(e) | Satisfăcătoare | Bună(e) | Foarte bună(e) | Nu e cazul |
|-----|--|-------------------|------------|----------------|---------|----------------|------------|
| D14 | Sănătatea dvs.   | 1                 | 2          | 3              | 4       | 5              | –          |
| D15 | Relațiile din familie  | 1                 | 2          | 3              | 4       | 5              | 98         |
| D16 | Locuința dvs.  | 1                 | 2          | 3              | 4       | 5              | –          |
| ... |  | 1                 | 2          | 3              | 4       | 5              |            |
| D27 | Posibilitățile existente de petrecere a timpului liber (de recreere) | 1                 | 2          | 3              | 4       | 5              | –          |

Sintaxa de etichetare a valorilor acestor variabile va fi :

**VALUE LABELS** d14 d15 d16 d17 d18 d19 d20 d21 d22 d23 d24 d25 d26 d27

1 „foarte proasta(e)”

2 „proasta(e)”

3 „satisfăcătoare”

4 „buna(e)”

5 „foarte buna(e)”

98 „nu e cazul”

Nu trebuie decât să notăm după **VALUE LABELS** lista de variabile care au aceleași variante de răspuns și aceleași etichete ale variantelor de răspuns. Simplu, nu ? Și mult mai rapid decât dacă am fi utilizat interfața grafică.

Ați observat, probabil, că nu folosesc diacritice în etichete. În această carte, în tabelele copiate din SPSS, am preferat să nu folosesc diacritice pentru a reproduce cât mai fidel senzația din timpul interacțiunii dvs. cu programul. Nu toți utilizatorii au computerele setate pentru a recunoaște diacriticele. De aceea, pentru a avea compatibilitate pe toate computerele, prefer să nu le utilizez. Cea mai neplăcută situație ar fi ca programul să nu le recunoască și să le înlocuiască cu un semn de întrebare sau cu un alt caracter. De asemenea, conform manualului programului, utilizarea diacriticelor poate crește considerabil dimensiunea bazei de date, ceea ce duce la creșterea timpului de deschidere a fișierului și de rulare a analizelor.

## 4.2. Definirea nonrăspunsurilor

Nonrăspunsurile (**missing values**) reprezintă absența răspunsului valid. Nonrăspunsurile pot fi clasificate în două tipuri generale: cele care țin de aplicarea chestionarului ca întreg persoanelor care ar trebui selectate conform schemei de eșantionare (*unit nonresponse*) și cele care țin de absența răspunsurilor la anumite întrebări din chestionar în cazul unei persoane selectate în eșantion (*item nonresponse*).

Primul tip de problemă apare, de exemplu, din cauza cadrelor de eșantionare care nu sunt actualizate sistematic, cum ar fi lista persoanelor cu drept de vot. Astfel, operatorul, când vizitează adresa primită, s-ar putea să nu mai găsească persoana inclusă în eșantion pentru că aceasta s-a mutat, a decedat etc. O altă cauză a nonrăspunsului de acest gen îl reprezintă dificultatea tot mai mare de a-i convinge pe oameni să răspundă la chestionare: aceștia nu au încredere în operatori, s-au plictisit din cauza solicitărilor frecvente primite de la diferite instituții care realizează astfel de cercetări, nu au încredere în modul cum sunt gestionate răspunsurile pe care le oferă etc. O analiză detaliată a acestor probleme este realizată de Ineke Stoop (2005) în lucrarea sa intitulată sugestiv *The Hunt for the Last Respondent*. Tot în această direcție a existat și există o preocupare constantă în diferite anchete comparative, cum ar fi *European Social Survey*, care oferă acces la o documentație vastă în această zonă și nu numai.

Al doilea tip de problemă apare, de exemplu, din cauza neatenției operatorului care sare peste o întrebare, refuzului de a răspunde al persoanei intervievate, care consideră întrebarea prea personală, modului cum a fost formulată o întrebare astfel încât respondentul care nu deține informația respectivă se vede nevoit să declare că nu știe răspunsul etc. În principiu, acest gen de nonrăspuns poate fi evitat prin modul cum sunt formulate întrebările și prin pregătirea riguroasă a operatorilor de teren. Însă, în realitate, multe chestionare conțin răspunsuri de tip „nu știu” sau „nu răspund”. Acestea nu sunt răspunsuri valide și trebuie tratate separat în baza de date. Și nonrăspunsul este însă un fel de răspuns, așa că, privind din perspectiva metodologului, ar fi util să realizăm un profil al acestor persoane pentru ca în cercetarea următoare să minimizăm aceste probleme.

În această secțiune, mă voi referi doar la al doilea tip de nonrăspuns (*item nonresponse*). Nu voi analiza problemele care îl generează, ci doar cum putem lucra în SPSS cu acest gen de date. În SPSS, ca și în alte programe de statistică, de altfel, nonrăspunsul este denumit *missing value*. De multe ori, în practică, am întâlnit mai frecvent denumirea în limba engleză, și nu cea în limba română. Acesta este doar rezultatul utilizării frecvente de către cercetători a programelor de analiză a datelor care au interfața în limba engleză. În bazele de date se folosesc coduri speciale pentru nonrăspunsuri. Cel mai adesea am întâlnit codurile 97, 98 și/ sau 99, respectiv derivate ale acestora: 7, 997, 9997, 8, 998, 9998, 9, 999, 9999

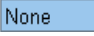
etc. În bazele de date internaționale se folosesc (și) alte coduri : -5, -4, -3, -2, -1. Nu trebuie să le folosiți numai pe acestea. Important este să utilizăm un cod pentru nonrăspuns care este foarte diferit de răspunsul valid. Să luăm câteva exemple :

- Variabila d4 din DCV 2010, „statutul ocupațional”, are patru răspunsuri valide : salariat (codul 1), pe cont propriu (codul 2), patron (codul 3) și zilier (codul 4). În mod normal, ar trebui să primim răspunsuri valide de la toți respondenții pentru că este o întrebare ușor de înțeles, cu variante clare. S-ar putea însă ca un respondent să nu dorească să declare statutul său ocupațional curent. Acest nonrăspuns va fi codificat cu 9, 99, 999 sau orice altă valoare similară sau putem utiliza codul -2, similar cercetării WVS 2012.
- Variabila d30 din DCV 2010, „Cum apreciați serviciul de pensii din România ? ”, are cinci răspunsuri valide : foarte prost (codul 1), prost (codul 2), satisfăcător (codul 3), bun (codul 4) și foarte bun (codul 5). Întrebarea este aplicată tuturor respondenților. Un respondent care nu are pensie sau nu cunoaște pe cineva care are pensie s-ar putea să declare că nu știe să evalueze acest sistem. Acest nonrăspuns va fi codificat cu 8, 98, 998 sau orice altă valoare similară sau putem utiliza codul -1 similar cercetării WVS 2012.
- Variabila d10 din DCV 2010, „starea civilă”, are șase răspunsuri valide : necăsătorit (nu a fost căsătorit niciodată) (codul 1), căsătorit (codul 2), divorțat (codul 3), separat (codul 4), văduv (codul 5) și altă situație (codul 6). Respondenții care aleg codurile 1, 3, 4, 5 sau 6 sunt rugați să răspundă la o întrebare suplimentară : „Aveți un partener de viață (cu care locuiți împreună, aveți menaj comun) ? ”. Cei care au răspuns codul 2, adică sunt căsătoriți, nu mai trebuie să răspundă la această întrebare. Nu li se aplică. Acesta este un tip aparte de nonrăspuns, denumit „nu e cazul”, care va fi codificat cu 7, 97, 997 sau orice altă valoare similară sau putem utiliza codul -3 similar cercetării WVS 2012.

În Access sau în programul pe care îl utilizăm pentru introducerea datelor, am definit deja aceste nonrăspunsuri pentru a ușura procesul de introducere a datelor și de curățare a bazei de date. Teoretic, nu ar trebui să mai introducem coduri în faza de curățare.

Unii cercetători preferă să nu instruiască SPSS că „nu știu” (98), „nu răspund” (99) sau „nu e cazul” (97) sunt nonrăspunsuri, lăsând acest lucru pentru momentul analizei pe care o va face. Alții preferă ca baza dată echipei de cercetare sau altor utilizatori să aibă deja nonrăspunsurile definite.

Nonrăspunsurile pot fi definite fie în **Variable View**, fie în sintaxă. Vă recomand a doua variantă.

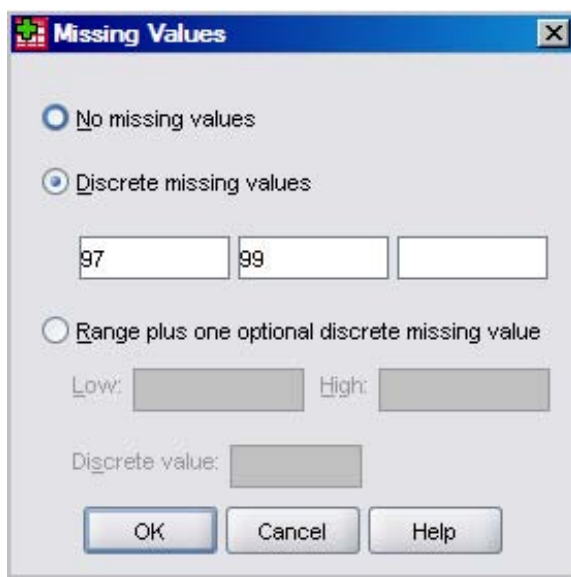
În **Variable View**, mergeți cu cursorul pe celula din dreptul variabilei care vă interesează (pe rând) și al coloanei **Missing** (pe coloană). Se vor activa cele trei puncte,  ..., pe care dăm click. După ce am realizat modificările, în locul



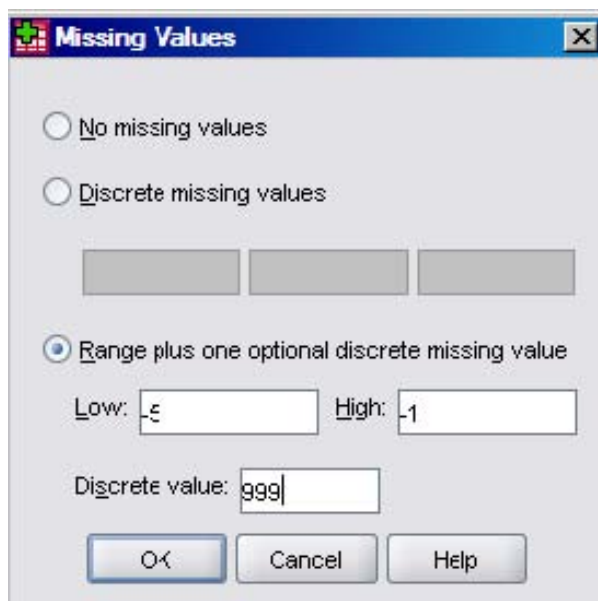
cuvântului **None**, vom observa valorile pe care le-am definit ca nonrăspunsuri. Se deschide fereastra din figura 4.5. În această fereastră, inițial, este selectată opțiunea **No missing values**.

**Figura 4.5.** Definirea nonrăspunsurilor în Variable View

(a)



(b)



Există două opțiuni pe care le putem bifa.

Opțiunea **Discrete missing values** este folosită atunci când variabila are maximum trei coduri de nonrăspuns, adică maximum trei tipuri de nonrăspuns : „nu știu”, „nu răspund”, „nu e cazul”. În figura 4.5a definim nonrăspunsurile pentru variabila d3, din DCV 2010, care conține informații despre anul de când respondentul are ocupația pe care a declarat-o la d2. Prin realizarea unui tabel de frecvență, am observat că d3 are două tipuri de nonrăspuns : „nu știu/nu răspund”, care a primit codul 99, și „nu e cazul”, care a primit codul 97. 463 de respondenți au primit codul 97 pentru că nu au o ocupație (sunt inactivi pe piața muncii), iar 165 nu au știut sau nu au dorit să precizeze anul de când au ocupația actuală.

Opțiunea **Range plus one optional discrete missing values** este folosită atunci când avem mai mult de trei tipuri de nonrăspuns : „nu știu”, „nu răspund”, „nu e cazul”, „întrebarea nu a fost adresată în anul respectiv” etc. De exemplu, în WVS 2012 avem coduri de la -5 la -1 : „missing : unknown” (codul -5), „not asked in survey” (-4), „not applicable” (-3), „no answer” (-2), „don’t know” (-1). La **Low** introducem -5, la **high** introducem -1. În figura 4.5b definim nonrăspunsurile pentru o variabilă care are coduri de nonrăspuns de la -5 la -1, dar și un cod 999. În principiu, această opțiune acoperă toate situațiile posibile.

De unde știm ce coduri trebuie să introducem în aceste celule? Am precizat deja că realizăm un tabel de frecvență pentru fiecare variabilă pentru care vrem să definim nonrăspunsurile. Să luăm ca exemplu variabila d15, „Cât de mulțumit sunteți de relațiile din familie?” (tabelul 4.1).

**Tabelul 4.1.** Tabel de frecvență înainte de definirea nonrăspunsurilor

| d15 relațiile din familie |                       |           |         |               |                    |
|---------------------------|-----------------------|-----------|---------|---------------|--------------------|
|                           |                       | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                     | 1 foarte proasta(e)   | 8         | .7      | .7            | .7                 |
|                           | 2 proasta(e)          | 14        | 1.2     | 1.2           | 1.9                |
|                           | 3 satisfacatoare      | 107       | 9.2     | 9.2           | 11.1               |
|                           | 4 buna(e)             | 666       | 57.4    | 57.4          | 68.5               |
|                           | 5 foarte buna(e)      | 289       | 24.9    | 24.9          | 93.4               |
|                           | 98 nu e cazul         | 42        | 3.6     | 3.6           | 97.0               |
|                           | 99 nu stiu/nu raspund | 35        | 3.0     | 3.0           | 100.0              |
|                           | Total                 | 1161      | 100.0   | 100.0         |                    |

Tabelul este realizat după ce am încheiat etapa de etichetare a variabilelor și a valorilor variabilelor. Observăm că această variabilă are cinci răspunsuri valide : „foarte proastă(e)” (codul 1), „proastă(e)” (codul 2), „satisfăcătoare” (codul 3), „bună(e)” (codul 4) și „foarte bună(e)” (codul 5). De asemenea, are două tipuri de nonrăspuns : „nu e cazul”, codul 98, și „nu știu/nu răspund”, cumulate în codul 99. Observați că nu există o regulă strictă care impune utilizarea acelorași coduri

în toate cercetările. Trebuie doar să existe o anumită consistență pentru a face mai ușoară tranziția de la o cercetare la alta. Cei care au răspuns „nu e cazul” s-au gândit, probabil, că întrebarea se referă la o relație maritală de tip soț-soție. Dacă inspectăm această ipoteză, observăm că toți cei 42 de respondenți din categoria „nu e cazul” (98) sunt necăsătoriți, divorțați sau văduvi. Așadar, ei au considerat că nu pot răspunde la această întrebare. În mod normal, în faza de curățare, dacă cercetătorului i se pare ciudat ca o persoană să nu răspundă la o întrebare, ar putea verifica chestionarele pentru o posibilă eroare de introducere sau ar putea chiar discuta cu operatorul de teren solicitând, uneori, refacerea chestionarului. Revenind la definirea nonrăspunsurilor, am aflat că trebuie să introducem în celulele **Discrete missing values** codurile 98 și 99.

Putem automatiza activitatea de definire a nonrăspunsurilor folosind sintaxa. Comanda este la fel de simplă ca celelalte două comenzi discutate: **VARIABLE LABELS** și **VALUE LABELS**. Comanda pentru nonrăspunsuri este **MISSING VALUES**. Mai exact, pentru variabila d15 comanda este:

**MISSING VALUES** d15 (98, 99).

La fel ca la **VALUE LABELS**, putem utiliza aceeași linie de comandă pentru mai multe variabile care au coduri similare la nonrăspunsuri. De exemplu, succesiunea de variabile d15-d27 se află în această situație. Așadar, comanda va arăta astfel:

**MISSING VALUES** d15 d16 d17 d18 d19 d20 d21 d22 d23 d24 d25 d26 d27 (98, 99).

Observăm cât de ușor putem defini nonrăspunsurile în acest meniu față de meniul **Variable View**, unde am fi dat mai multe clickuri pentru fiecare variabilă în parte. În plus, oricând dorim, putem consulta sintaxa, reamintindu-ne ce am lucrat sau pentru a o rula din nou pe o bază „curată”.

Pentru că orice proces de învățare presupune căutare de informație, vă invit să aflați ce element din comanda **MISSING VALUES** puteți șterge fără a afecta rezultatul final.

### 4.3. Verificarea introducerii eronate a unor coduri

Dacă am folosit un program de introducere a datelor care restricționează operatorul de introducere să introducă greșit o valoare în afara amplitudinii răspunsurilor posibile, în principiu, atunci putem sări această etapă, deși niciodată nu strică o verificare.

Verificarea este o operație simplă care presupune doar inspectarea tabelelor de frecvență pentru fiecare variabilă din baza de date. Deja am precizat că aceste tabele se realizează din meniul **Analyze > Descriptive statistics > Frequencies**.

Trebuie să comparăm ce apare în tabel cu ceea ce este scris în chestionar. Această operațiune poate fi realizată fie înainte de definirea nonrăspunsurilor, fie ulterior. Este util, în schimb, să fie încheiată operațiunea de etichetare a valorilor variabilelor pentru a vedea ce reprezintă fiecare cod.

#### 4.4. Validarea logică prin urmărirea filtrelor și a unor întrebări factuale

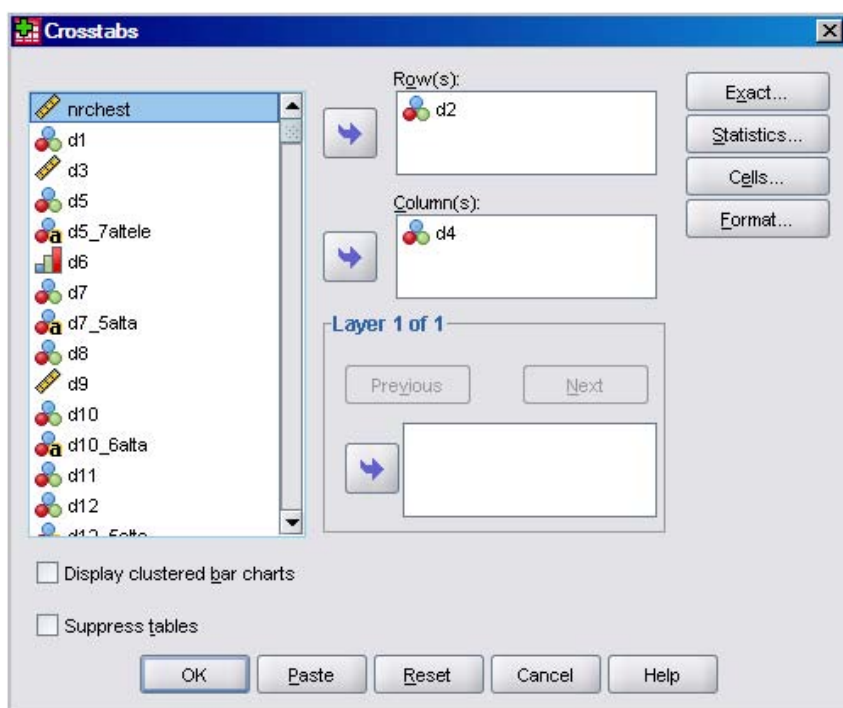
Cei care au declarat că sunt elevi (codul 6), pensionari (codul 7), casnici (codul 8) sau șomeri (codul 9) la variabila d2, ocupația principală actuală, trebuiau să răspundă apoi direct la întrebarea d6. Așadar, aceștia trebuie să aibă, la întrebările dintre d2 și d6, un cod de nonrăspuns de tipul „nu e cazul”. Pur și simplu, întrebările dintre d2 și d6 nu li se aplică. Dacă am lucrat corect în programul de introducere a datelor sau în alt program similar nu ar trebui să apară erori de introducere. Putem verifica simplu dacă filtrul a fost respectat, realizând un tabel de contingență folosind meniul **Analyze > Descriptive statistics > Crosstabs** între d2 și fiecare dintre întrebările de până la d6. În tabelul 4.2 este prezentat un exemplu de încrucișare între d2 (ocupația principală actuală) și d3 (statutul ocupațional). Conform chestionarului, în celulele rezultate din intersecția dintre rândurile ce conțin codurile 6, 7, 8 și 9 și coloanele date de răspunsurile valide la d4 și codul 99 („nu știu/nu răspund”) ar trebui să apară valoarea 0, adică nici o persoană. Observăm că aici filtrul este respectat: apar persoane doar la intersecția dintre codurile 6-9 la d2 și codul 97 („nu e cazul”) la d4.

**Tabelul 4.2.** Tabel de contingență ce verifică un filtru, dar este folosit pentru validare logică (1)

| <b>d2 ocupatia dvs. actuala (principala) * d4 statutul ocupational Crosstabulation</b> |                                     |                         |                      |             |          |     |    |       |
|--|-------------------------------------|-------------------------|----------------------|-------------|----------|-----|----|-------|
| Count  |                                     |                         |                      |             |          |     |    |       |
|  |                                     | d4 statutul ocupational |                      |             |          |     |    | Total |
|  |                                     | 1<br>salariat           | 2 pe cont<br>propriu | 3<br>patron | 4 zilier | 97  | 99 |       |
| d2<br>ocupatia<br>dvs.<br>actuala<br>(princi-<br>pala)                                 | 1 agricultor                        | 2                       | 228                  | 0           | 29       | 0   | 0  | 259   |
|  | 2 muncitor (meserias)               | 215                     | 14                   | 1           | 12       | 0   | 5  | 247   |
|  | 3 tehnician, maistru,<br>functionar | 72                      | 2                    | 0           | 0        | 0   | 0  | 74    |
|  | 4 ocupatie cu studii<br>superioare  | 98                      | 8                    | 0           | 0        | 0   | 0  | 106   |
|  | 6 elev, student                     | 0                       | 0                    | 0           | 0        | 52  | 0  | 52    |
|  | 7 pensionar                         | 0                       | 0                    | 0           | 0        | 267 | 0  | 267   |
|  | 8 casnic                            | 0                       | 0                    | 0           | 0        | 62  | 0  | 62    |
|  | 9 acum sunt somer                   | 0                       | 0                    | 0           | 0        | 82  | 0  | 82    |
|  | 10 patron                           | 0                       | 0                    | 12          | 0        | 0   | 0  | 12    |
| Total  |                                     | 387                     | 252                  | 13          | 41       | 463 | 5  | 1161  |

Revenim la modul în care am realizat tabelul de contingență. Deși acestui subiect i se dedică o secțiune specială, cred că este util să vedem pașii elementari în realizarea acestui tip de tabel și aici. Accesând meniul **Analyze > Descriptive statistics > Crosstabs** se deschide fereastra din figura 4.6. Un tabel de contingență are două variabile. O variabilă, prin categoriile ei, dă rândurile tabelului. Cealaltă variabilă, prin categoriile ei, dă coloanele tabelului. Fiecare celulă din tabel ne arată numărul persoanelor care se regăsesc în două categorii simultan : 2 persoane sunt agricultori salariați, 215 sunt muncitori salariați, 8 persoane au o ocupație care necesită studii superioare și lucrează pe cont propriu etc. Care celulă din tabel prezintă o informație inconsistentă ? Cum puteți explica această inconsistență și ce ar trebui să faceți pentru a o corecta ?

**Figura 4.6.** Crosstabs : realizarea unui tabel de contingență



Prefer să introduc în rând (celula **Row**) variabila care are cele mai multe variante de răspuns pentru a rezulta un tabel pe verticală, ușor de inserat într-o pagină A4 orientată portret. Pe coloană (celula **Column**) introduc cealaltă variabilă. Pentru ce avem nevoie acum este suficient să apăsăm **OK**. Va rezulta tabelul cu frecvențele absolute, adică cu celulele arătând numărul de agricultori care sunt salariați, numărul de agricultori care lucrează pe cont propriu, numărul de muncitori care sunt salariați ș.a.m.d.

Ați reușit să observați inconsistența ? O persoană care a declarat că ocupația sa este muncitor a indicat ca statut ocupațional faptul că este patron. Aici avem de-a face cu o validare logică a informației conținute de baza de date. Cei care au ales varianta 2 la d2 puteau răspunde la d4, deci nu avem un filtru. În schimb, logic ar fi ca un muncitor să nu se declare patron. Primul lucru pe care ar trebui să îl facem este să consultăm chestionarul completat. Dacă nu găsim răspunsul la această inconsistență în el, atunci va trebui să discutăm cu operatorul de teren pentru a ne clarifica situația.

Să presupunem că am identificat de unde vine această problemă. Am aflat că, de fapt, persoana respectivă are ocupația de patron. Deci a fost introdus greșit în baza de date codul 2 în loc de codul 10 la d2. Așadar trebuie să modificăm informația în baza de date. Acest lucru se face prin recodificarea variabilei d2. Dar pentru a face recodificarea trebuie să aflăm id-ul unic al celui respondent pentru a fi siguri că modificăm numai ce ne interesează. Acest lucru se face urmând mai mulți pași :

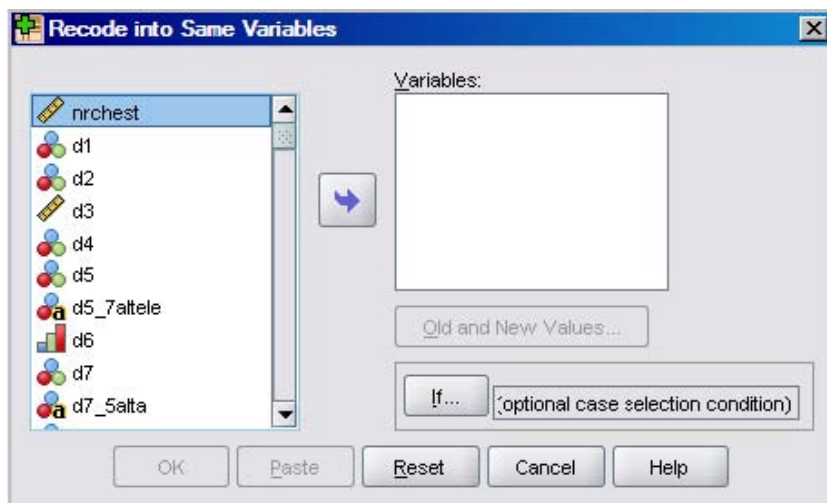
- filtrăm baza de date astfel încât să rămână activ doar cazul care are codul 2 la d2 și codul 3 la d4, adică muncitorul care a declarat statutul patron. Filtrul este : **d2 = 2 & d4 = 3**. Mergem în **Data > Select Cases > If condition is satisfied > If > introducem condiția > Continue > OK** ;
- realizăm un tabel de frecvență pentru variabilele care sunt folosite pentru condiție și pentru variabila care conține id-ul unic, aici nrchest ;
- verificăm dacă filtrul activ este cel dorit ;
- citim tabelul de frecvență pentru variabila nrchest și aflăm că acel caz are id-ul 312 ;
- dezactivăm filtrul.

Acum putem trece la recodificarea variabilei d2, pentru că ea conține eroarea. Recodificarea se va face din meniul **Transform > Recode into Same Variables**. Acest meniu va înlocui, pentru cazul cu id-ul 312, codul 2 cu codul 10. Figura 4.7 prezintă etapele acestui proces :

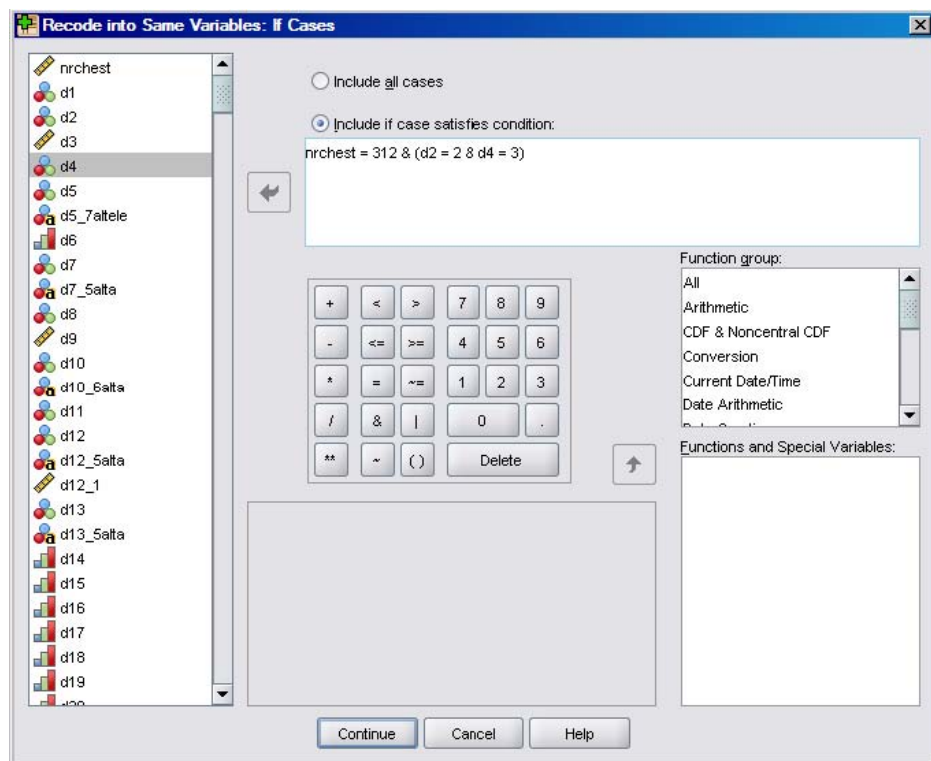
- Apăsând butonul **If > Include if case satisfies condition**, punem condiția ca modificarea să fie realizată doar pentru cazul cu id-ul 312. Aici am notat și că **d2 = 2 & d4 = 3**. Apăsăm **Continue** (figura 4.7b).
- Apăsăm butonul **Old and New Values**.
- Înlocuim codul 2 (**Old Value > Value**) cu codul 10 (**New Value > Value**). Apăsăm butonul **Add**, apoi apăsăm butonul **Continue** (figura 4.7c).
- Am revenit în fereastra inițială, unde apăsăm butonul **OK**.
- Refacem tabelul de contingență dintre d2 și d4, pentru a verifica dacă modificarea s-a făcut conform așteptărilor (tabelul 4.3).

**Figura 4.7. Recode into Same Variables**

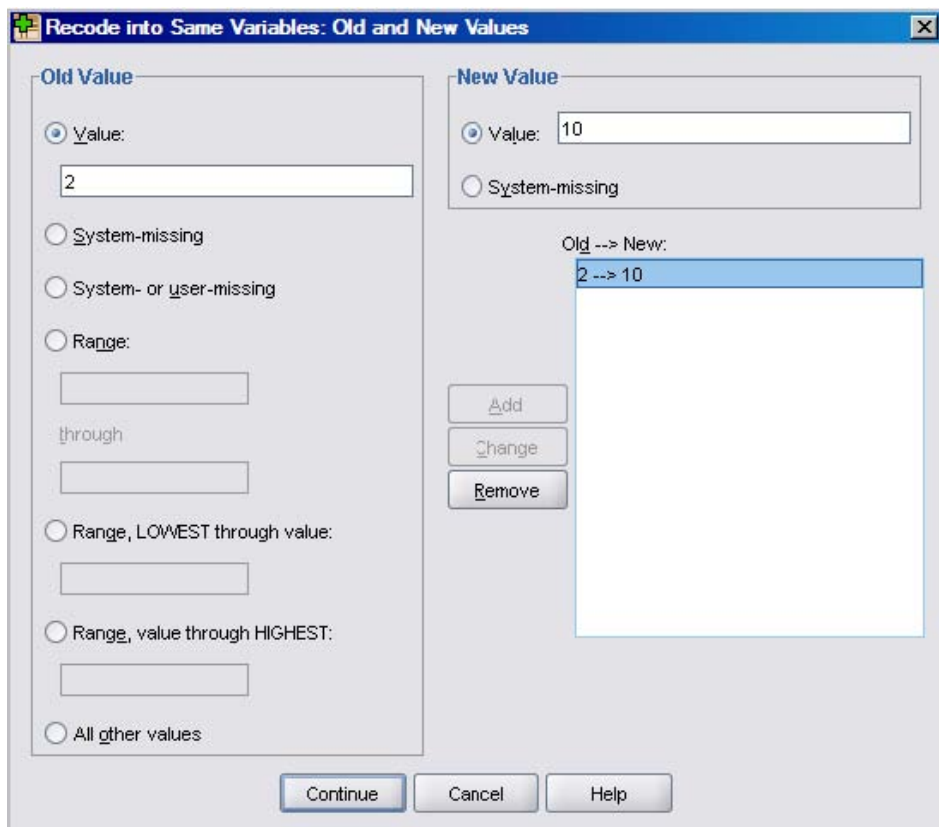
(a)



(b)



(c)



Utilizarea condițiilor în meniul **Recode into Same Variables** nu este obligatorie. Acest lucru a fost impus de situația discutată. De multe ori însă folosim doar comenzile activate de butonul **Old and New Values**. Fereastra care se deschide are mai multe secțiuni (figura 4.7b). În secțiunea **Old Value** introduc valorile inițiale: cele pe care dorim să le recodificăm. În secțiunea **New Value** introduc valorile noi: cele în care vor fi recodificate valorile inițiale. În secțiunea **Old --> New**, după ce apăsăm butonul **Add**, apar modificările pe care dorim să le facem. Aici am dorit să modificăm doar un cod: 2 în 10. De aceea am folosit **Old Value > Value**. Dacă am fi vrut să modificăm o serie de numere consecutive, să zicem 2-6 în 10, atunci am fi folosit **Old Value > Range 2 through 6**. Dacă am fi vrut să modificăm o serie de numere consecutive de la cea mai mică valoare la o valoare anume, să zicem de la 1 (valoarea minimă) la 4 (valoarea specifică), atunci am fi folosit **Old Value > Range, LOWEST through value: 4**. Dacă am fi vrut să modificăm o serie de numere consecutive de la o valoare specifică la cea mai mare din serie, să zicem de la 4 (valoarea specifică) la 10 (valoarea maximă), atunci am fi folosit **Old Value > Range, value through**



**HIGHEST: 10.** Vom discuta despre recodificare în altă secțiune a lucrării. Ce trebuie reținut aici este că folosirea meniului **Recode into Same Variables** suprascrie informația inițială. Deci atenție la ce modificări doriți să faceți.

**Tabelul 4.3.** Tabel de contingență care verifică un filtru, dar este folosit și pentru validare logică (2)

| d2 ocupatia dvs. actuala (principala) * d4 statutul ocupational Crosstabulation |  |                         |                         |             |             |     |    |       |
|---|--|-------------------------|-------------------------|-------------|-------------|-----|----|-------|
| Count   |  |                         |                         |             |             |     |    |       |
|   |  | d4 statutul ocupational |                         |             |             |     |    | Total |
|   |  | 1<br>salariat           | 2 pe<br>cont<br>propriu | 3<br>patron | 4<br>zilier | 97  | 99 |       |
| d2 ocupatia<br>dvs. actuala<br>(principala)                                     | 1 agricultor                           | 2                       | 228                     | 0           | 29          | 0   | 0  | 259   |
|   | 2 muncitor<br>(meserias)               | 215                     | 14                      | 0           | 12          | 0   | 5  | 246   |
|   | 3 tehnician,<br>maistru,<br>functionar | 72                      | 2                       | 0           | 0           | 0   | 0  | 74    |
|   | 4 ocupatie cu<br>studii superioare     | 98                      | 8                       | 0           | 0           | 0   | 0  | 106   |
|   | 6 elev, student                        | 0                       | 0                       | 0           | 0           | 52  | 0  | 52    |
|   | 7 pensionar                            | 0                       | 0                       | 0           | 0           | 267 | 0  | 267   |
|   | 8 casnic                               | 0                       | 0                       | 0           | 0           | 62  | 0  | 62    |
|   | 9 acum sunt<br>somer                   | 0                       | 0                       | 0           | 0           | 82  | 0  | 82    |
|   | 10 patron                              | 0                       | 0                       | 13          | 0           | 0   | 0  | 13    |
| Total   |  | 387                     | 252                     | 13          | 41          | 463 | 5  | 1161  |

Așadar, pe lângă verificarea filtrelor, realizăm și validarea logică prin încrucișarea unor variabile factuale. O variabilă factuală culege informații concrete care nu țin de valori, atitudini, opinii, credințe, evaluări. De exemplu, sexul sau vârsta respondentului sunt variabile factuale. Tot variabile factuale sunt salariul măsurat într-o unitate monetară, suprafața locuinței în metri pătrați, numărul de copii etc. În chestionare, din cauza neatenției operatorului sau poate chiar dintr-o scăpare de design a cercetătorului, se mai întâmplă ca o persoană să declare ceva la o variabilă factuală, acel ceva fiind incompatibil cu ce declară la altă variabilă factuală aflată într-o relație logică cu cea dintâi. Un bărbat nu are voie să răspundă la întrebarea „Ați făcut vreodată avort?”. În schimb, are voie să răspundă la întrebarea „Partenera dvs. de viață a făcut vreodată avort?”. Cel care a declarat că nu suferă de vreo boală nu are voie să răspundă la întrebarea „Suferiți de o boală cronică?”. Dincolo de cele două tipuri de erori enunțate mai există o situație care, într-un fel, ține de designul chestionarului, deci este o problemă a cercetătorului. Realitatea din teren s-ar putea să fie mai complexă decât cea pe care o

cunoaște sau și-o imaginează cercetătorul. De exemplu, un cercetător s-ar putea aștepta ca o persoană care declară că este pensionar să nu mai ofere un răspuns valid la rubrica „Vă rugăm să ne spuneți care a fost suma de bani încasată luna trecută din salarii”, ci doar la rubrica „Vă rugăm să ne spuneți care a fost suma de bani încasată luna trecută din pensii”. Salariul, teoretic, este specific unei persoane angajate formal, cu contract de muncă. Totuși, salariul poate fi atribuit și persoanelor care nu sunt angajate formal, ci prestează diferite servicii informal („la negru”). Când ne gândim la un pensionar ne imaginăm că salariul acestuia este pensia, deci nu mai prestează servicii, cel puțin formalizate. Nivelul de trai redus din România și, implicit, al pensiilor îi determină pe mulți pensionari să lucreze informal. De exemplu, un pensionar se poate „angaja” ca paznic de noapte la o firmă. Acesta primește o pensie, dar și un salariu, chiar dacă acel salariu nu este înregistrat legal. Aici intervine altă problemă : să presupunem că cercetătorul admite că acest gen de situații este veridic, astfel încât îi va adresa întrebarea referitoare la salariu și pensionarului. Pensionarul, în schimb, fiind conștient că salariul său nu este înregistrat legal, s-ar putea să refuze să răspundă la întrebarea legată de salariu și să accepte să răspundă doar la întrebarea legată de pensie. Astfel apare nonrăspunsul și, implicit, discuția se mută în zona de distorsiune a realității, de modificare a reprezentativității eșantionului.

Revenind la problema validării logice prin încrucișarea variabilelor factuale, primul pas ce trebuie făcut este să identificăm în chestionar toate interacțiunile posibile dintre variabilele factuale. Apoi, realizând tabele de contingentă, așa cum am discutat deja, scanăm datele pentru erori. Termenul „eroare” este poate prea tranșant. Cercetătorii trebuie să consulte chestionarele tipărite și, eventual, să contacteze din nou respondentul pentru clarificări. Abia apoi se intervine în baza de date. Validarea logică poate fi inclusă chiar în partea de design a cercetării și chestionarului. De exemplu, într-un studiu prin care se dorea estimarea incidenței consumului diferitelor tipuri de droguri, cercetătorul a introdus în lista de droguri și unul fictiv. Dacă în chestionar apăreau răspunsuri valide la acest drog (respondentul „spunea” că a consumat, cu o anumită frecvență, în anumite condiții etc.), atunci cercetătorul afla imediat că operatorul de teren nu a fost onest și a completat el însuși acel chestionar.

În DCV 2010, există variabila d39 : „Caracterizați măsura în care puteți influența luarea deciziilor în organizația în care lucrați” cu variantele de răspuns „foarte scăzută” (codul 1), „scăzută” (codul 2), „satisfăcătoare” (codul 3), „ridicăată” (codul 4), „foarte ridicată” (codul 5). În chestionar, în dreptul acestei variabile, există și varianta „nu e cazul” (codul 98). Logica este simplă : o persoană care nu are un loc de muncă nu poate să evalueze măsura în care are libertate de decizie acolo (organizație este un termen generic pentru toate locurile de muncă fie că acestea sunt în firme, instituții publice, ONG-uri etc.). Prin urmare, trebuie să verificăm legătura logică cu variabila factuală d2, „ocupația dvs actuală principală”, care are 10 variante de răspuns : „agricultor” (codul 1), „muncitor (meseriaș)” (codul 2), „tehnician, maistru, funcționar” (codul 3),

„ocupație cu studii superioare” (codul 4), „altă ocupație” (codul 5), „elev, student” (codul 6), „pensionar” (codul 7), „casnic” (codul 8), „acum sunt șomer” (codul 9) și „patron” (codul 10). Logic ar fi ca elevii/studenții, pensionarii, casnicii și șomerii, adică codurile 6, 7, 8 și 9, să nu aibă răspunsuri valide la d39. Tabelul de contingență de mai jos (tabelul 4.4) ne arată o încălcare a acestei logici : există 2 studenți care evaluează libertatea de decizie ca fiind ridicată, 3 șomeri care o evaluează ca fiind foarte scăzută și 4 șomeri care o evaluează ca fiind scăzută, 1 pensionar care o evaluează ca fiind scăzută.

**Tabelul 4.4.** Validare logică : tabel de contingență

| <b>d2 ocupatia dvs. actuala (principala) * d39 masura in care puteti influenta luarea deciziilor in organizatia in care lucrați Crosstabulation</b> |  |  |                   |                          |                    |                              |                     |    |       |
|---|--|--|-------------------|--------------------------|--------------------|------------------------------|---------------------|----|-------|
| Count   |  |  |                   |                          |                    |                              |                     |    |       |
|   |  | d39 masura in care puteti influenta luarea deciziilor in organizatia in care lucrați |                   |                          |                    |                              |                     |    | Total |
|   |  | 1<br>foarte<br>sca-<br>zuta  | 2<br>sca-<br>zuta | 3<br>satisfa-<br>catoare | 4<br>ridi-<br>cata | 5<br>foarte<br>ridi-<br>cata | 98<br>nu e<br>cazul | 99 |       |
| d2 ocupatia<br>dvs. actuala<br>(principala)   | 1 agricultor                           | 8  | 15                | 26                       | 4                  | 1                            | 197                 | 8  | 259   |
|   | 2 muncitor<br>(meserias)               | 43   | 66                | 77                       | 29                 | 9                            | 18                  | 4  | 246   |
|   | 3 tehnician,<br>maistru,<br>functionar | 11   | 18                | 22                       | 18                 | 2                            | 3                   | 0  | 74    |
|   | 4 ocupatii<br>cu studii<br>superioare  | 9  | 21                | 31                       | 23                 | 11                           | 11                  | 0  | 106   |
|   | 6 elev,<br>student                     | 0  | 0                 | 0                        | 2                  | 0                            | 50                  | 0  | 52    |
|   | 7 pensionar                            | 0  | 1                 | 0                        | 0                  | 0                            | 266                 | 0  | 267   |
|   | 8 casnic                               | 0  | 0                 | 0                        | 0                  | 0                            | 62                  | 0  | 62    |
|   | 9 acum sunt<br>somer                   | 3  | 4                 | 0                        | 0                  | 0                            | 73                  | 2  | 82    |
|   | 10 patron                              | 3  | 0                 | 1                        | 1                  | 2                            | 5                   | 1  | 13    |
| Total   |  | 77   | 125               | 157                      | 77                 | 25                           | 685                 | 15 | 1161  |

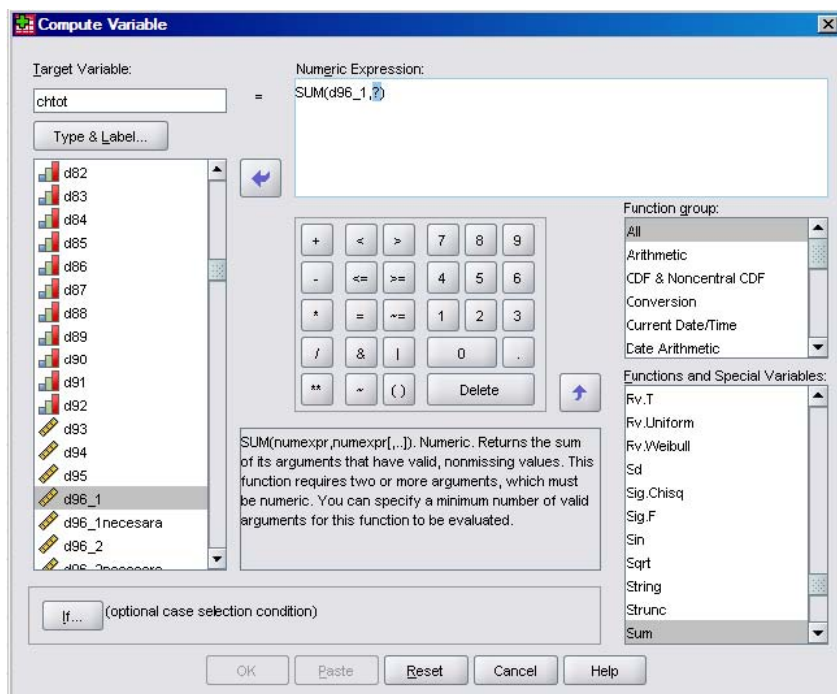
Încă o dată : modificările nu se fac automat. Am verificat chestionarele și am constatat următoarele situații : (1) studenții sunt, de fapt, persoane care au ocupații cu studii superioare (să nu uităm că d2 cere răspuns unic) ; (2) cei 3 șomeri sunt agricultori, iar (3) cei 4 șomeri sunt muncitori (meseriași). Așadar, trebuie să recodificăm variabila d2 pentru cazurile acestea. Dacă dorim să procedăm ca mai devreme, care sunt pașii pe care trebuie să îi parcurgeți pentru realizarea modificărilor ?

## 4.5. Construirea unor variabile noi

Aceasta poate fi sau nu o parte a procesului de curățare. Unii cercetători preferă să aibă anumite variabile create încă de la început, alții consideră că le pot construi singuri pe măsură ce este nevoie de ele. În DCV 2010, există întrebarea d96 prin care respondenții sunt rugați să raporteze, în lei, sumele cheltuite pentru 9 variabile: alimente; băuturi alcoolice, țigări; îmbrăcăminte, încălțăminte; pentru locuință: chirie, întreținere, reparații, abonamente, rate; transport; îngrijirea sănătății; cultură, școală, cărți, spectacole; pensie alimentară; alte cheltuieli. În unele analize, am putea fi interesați să lucrăm cu variabila care conține informații despre cheltuielile totale ale gospodăriei în luna precedentă. Aceasta este, evident, suma tuturor acestor nouă variabile. Suma aceasta va deveni o nouă variabilă în baza de date. Fiecare respondent va avea, în dreptul său, valoarea însumată a tuturor cheltuielilor efectuate.

Pentru a calcula această sumă, folosim meniul **Transform > Compute** (figura 4.8). Decidem că numele variabilei va fi „chtot”. Prefer numele scurte, formate doar din litere și, eventual, cifre, pentru că folosesc anumite programe de statistică specializate cum ar fi HLM (Raudenbush *et al.*, 2011), care solicită aceste specificații. Dacă nu le respect, programul va redenumi variabilele și, de multe ori, această operație automată creează nume cu care este greu de lucrat.

Figura 4.8. Crearea de noi variabile (Compute)



La **Target Variable** vom scrie numele variabilei pe care o realizăm, aici chto.

La **Numeric Expression** vom scrie formula care ne dă nouă variabilă. Putem folosi două abordări, în funcție de necesități: (1) folosim funcțiile pe care le oferă SPSS, cum ar fi funcția **SUM()**, sau (2) introducem noi o expresie de tipul  $a+b+c+\dots+n$ . Pentru moment, folosim funcția **SUM()** pe care o aducem în **Numeric Expression** din secțiunea **Functions and Special Variables**. Mai întâi, dăm click în secțiunea **Function group** pe **All**, pentru a se activa funcțiile din secțiunea **Functions and Special Variables**. Apoi căutăm funcția **SUM()** folosind același procedeu ca în lista de variabile din orice meniu ori, pur și simplu, utilizând scroll-ul. Când o găsim, dăm dublu click pe ea și vom vedea că va apărea în **Numeric Expression**. Inițial, ea arată astfel: **SUM(?, ?)**. Ștergem semnele de întrebare, căutăm în lista de variabile din stânga ceea ce ne interesează, aici succesiunea **d96\_1-d96\_9**, dăm, pe rând, dublu click pe variabile sau le introducem cu săgeata în dreapta, punând virgulă între ele. Apăsăm **OK**. Putem merge în **Data View** sau **Variable View** să vedem variabila. Variabilele noi sunt create la sfârșitul bazei de date. Putem, în loc de funcția **SUM()**, să adunăm pur și simplu variabilele respective. Adică, în **Numeric Expression**, să fi scris: **d96\_1 + d96\_2 + d96\_3 + d96\_4 + d96\_5 + d96\_6 + d96\_7 + d96\_8 + d96\_9**. Rezultatul este diferit și vine din modul în care SPSS tratează nonrăspunsurile. Să presupunem că la aceste întrebări există persoane care au refuzat să răspundă sau au declarat că nu știu să răspundă. Aceste valori nu sunt valide și nu vor fi luate în calcul la sumă dacă le-am definit în coloana **Missing** din **Variable View** sau folosind comanda **MISSING VALUES** în sintaxă. Dacă folosim funcția **SUM()**, atunci variabila chto va conține suma variabilelor chiar dacă, la una sau mai multe dintre ele, respondentul nu a indicat o valoare validă, ci a oferit, în schimb, un nonrăspuns. Dacă folosim adunarea, atunci variabila chto va conține suma pentru respondenții care au oferit răspunsuri valide la toate variabilele din șir, ștergându-i pe ceilalți. Deci în ultima variantă o să avem mai puține valori valide în variabila nou-creată, pentru că este folosită doar informația completă, pe când în prima variantă o să avem mai multe valori valide pentru că este folosită toată informația disponibilă. Care variantă este corectă? Răspunsul nu este atât de evident. Aici, la cheltuieli, am putea folosi și informația parțială notând totuși în lucrarea pe care o scriem această limită a analizei. Dacă avem o variabilă latentă, cum ar fi atitudinea față de fumat măsurată printr-o scală compusă care conține 5 itemi, atunci poate ar fi bine să folosim informația completă: altfel, scorul, adică atitudinea, s-ar baza pe o măsurătoare incompletă. Cel mai onest ar fi să construim ambele variabile, să rulăm analizele dorite cu ambele variabile separat și să vedem dacă rezultatele se schimbă substanțial.

Variabilele nou-create de noi nu sunt etichetate automat. Va trebui să rulăm de fiecare dată sintaxele **VARIABLE LABELS** și **VALUE LABELS**, în funcție de nevoi.

O altă situație, destul de frecvent întâlnită, în care se impune folosirea meniului **Transform > Compute** se referă la calcularea vârstei. În chestionar, respondentul nu este întrebat ce vârstă are, ci în ce an s-a născut. În analize însă, suntem

interesați să lucrăm cu vârsta, de aceea va trebui să creăm această variabilă în baza de date. Vârsta va fi egală cu anul aplicării chestionarului minus anul nașterii. În DCV 2010, respondentul este rugat să își declare vârsta în ani împliniți, situație care nu se aplică aici. Există însă variabila d3 care înregistrează anul din care respondentul are ocupația declarată la d2. În analize ne interesează să lucrăm cu variabila vechime în muncă măsurată în ani. Decidem să creăm această variabilă care se va numi „vechime”. Mai întâi, trebuie să ne asigurăm că la d3 sunt definite nonrăspunsurile. Observăm că în această bază avem codul 97, aplicat celor care nu au o ocupație în prezent, și codul 99, aplicat celor care nu au vrut să răspundă sau nu au știut unde să se încadreze în variantele puse la dispoziție de cercetător. Dacă nu facem acest lucru, vor fi luate în considerare la calcul și aceste coduri, noua variabilă conținând informații eronate. Evident va trebui să o etichetăm pentru a ști în continuare ce reprezintă. În cadrul unei analize, lucrăm cu multe variabile și este foarte ușor să uităm ce am făcut anterior, mai ales dacă lăsăm o pauză de câteva zile între început și sfârșit și lucrăm în mai multe proiecte simultan.

Posibilitățile pe care ni le oferă comanda **COMPUTE** sunt numeroase. Fiecare le va folosi pe cele care îi sunt utile în analize.

## 4.6. Exerciții

Pentru aceste exerciții utilizăm baza de date și/sau chestionarul World Values Survey 2012 rezultat(ă/e) în urma aplicării chestionarului în România. Baza de date poate fi descărcată de pe pagina de internet a *Grupului Românesc pentru Studiul Valorilor Sociale* (<http://www.romanianvalues.ro>).

1. Deschideți baza de date finală creată la exercițiul 13 din capitolul 2. Definiți proprietățile variabilelor din baza de date în Variable View.
2. Realizați câte un tabel de frecvență pentru fiecare variabilă din baza de date. Există coduri introduse eronat? Dacă da, cum explicați această greșală?
3. Identificați în cele patru pagini de chestionar alese în exercițiile din capitolul 2 întrebările filtru. Verificați dacă filtrele au fost respectate.
4. Identificați, în cele patru pagini de chestionar, variabile care pot fi folosite în procesul de validare logică. Verificați dacă există situații în care logica a fost încălcată.
5. Deschideți baza de date WVS 2012. Realizați câte un tabel de frecvență pentru fiecare variabilă din chestionar. Există coduri introduse eronat? Dacă da, cum explicați această greșală?
6. Identificați în chestionarul WVS 2012 întrebările filtru. Verificați dacă filtrele au fost respectate.
7. Identificați în chestionarul WVS 2012 variabile care pot fi folosite în procesul de validare logică. Verificați dacă există situații în care logica a fost încălcată.

## 5. Gestionarea variabilelor

Despre acest subiect am mai discutat. Am învățat să modificăm o variabilă folosind **Transform > Recode into Same Variables**. De asemenea, am învățat să creăm o variabilă nouă folosind o funcție sau o formulă, prin intermediul **Transform > Compute**. În acest capitol dezvoltăm acest subiect. O mare parte din activitatea de analiză cantitativă a datelor este destinată pregătirii variabilelor.

Voicu, Rusu și Comșa (2013) vor să explice solidaritatea românilor. Solidaritatea este, pentru aceștia, o atitudine față de alte persoane care denotă cooperare, interes, preocupare, sprijin etc. Solidaritatea este măsurată printr-un scor compozit obținut prin cumularea răspunsurilor la mai mulți itemi. Factorii care determină solidaritatea sunt orientarea postmaterialistă sau materialistă, religiozitatea și comportamentul religios, identificarea națională, încrederea generalizată, individualismul, clasa socială, vârsta, venitul, educația, sexul și tipul de localitate de rezidență. Analiza prin care doresc să testeze ipotezele este regresia liniară multiplă, tehnică prezentată într-un capitol al acestei cărți. Observăm că modelul explicativ propus de autori este destul de complex. Fiecare variabilă din model, începând cu cea dependentă (solidaritatea), trebuie pregătită pentru analiză. Pregătirea se va face ținând cont și de caracteristicile pe care le pot avea variabilele în analiza de regresie liniară. Din acest motiv, informațiile despre cum gestionăm variabilele capătă un rol esențial în procesul de analiză cantitativă.

În acest capitol vom discuta despre meniul **Transform > Recode into Different Variables** și vom afla câteva informații noi despre meniul **Transform > Compute**.

### 5.1. Crearea unei alte variabile utilizând meniul Recode into Different Variable

Înainte de a începe analiza datelor, vă recomand să salvați într-un loc sigur baza de date în forma primită de la cei care au curățat-o. Aceasta va fi baza de referință la care apelăm atunci când am pierdut informații din copia pe care lucrăm.

De exemplu, eu prefer să șterg codurile de nonrăspuns din baza de date, lăsând celulele respective goale. Astfel SPSS le va trata tot timpul ca **missing values** : nu mai există riscul să obțin rezultate greșite, pentru că am uitat să le definesc



în coloana **Missing** din **Variable View** sau folosind comanda **MISSING VALUES** din sintaxă. Această preferință poate fi satisfăcută folosind meniul **Transform > Recode into Same Variables**. Dar odată rulată comanda, am șters acea informație din variabilele cu care lucrez. Dacă, ulterior, doresc să realizez un profil al celor care au declarat că nu știu răspunsul la o întrebare și să îl compar cu profilul celor care nu vor să răspundă la aceeași întrebare, atunci nu mai pot face acest lucru. Să ne gândim la toate veniturile unei persoane. Aceasta poate să câștige bani din salariul la principalul loc de muncă, dar poate avea și un loc de muncă secundar, unde lucrează pe proiect. De asemenea, poate avea un cont de economii și astfel primește lunar o dobândă. Toate aceste venituri se adună și rezultă venitul lunar total al persoanei respective. Cercetătorul dorește să estimeze venitul mediu al românilor pentru luna februarie a anului 2014. Va pune în chestionar o rubrică de forma : „Dacă adunați veniturile din toate sursele, vă rog să îmi spuneți câți bani ați câștigat dvs. personal în luna februarie”. Respondentul trebuie să aproximeze o sumă dacă nu o cunoaște pe cea exactă. În teren, unii respondenți ne oferă un răspuns. Alții, în schimb, refuză să facă acest lucru. Motivele sunt multiple : lucrează „la negru”, operatorul nu le inspiră încredere să declare o informație atât de personală etc. În fine, cei care sunt mai puțin preocupați de gestionarea lunară a veniturilor lor s-ar putea să nu știe și, decât să ofere un răspuns greșit, preferă să aleagă această variantă de răspuns. În baza de date, la introducere, fiecare dintre aceste situații primește codul corespunzător. Deci variabila va avea valori de la 0 la cel mai mare venit și codurile 98 (nu știu) și 99 (nu răspund). Ipoteza mea este că cei care au declarat că nu știu au alte caracteristici decât cei care au refuzat să răspundă. Dacă vreau să le compar caracteristicile și am șters codurile, fără să fi păstrat o copie a bazei originale, nu mai pot face acest lucru. Concluzia : poate ar fi fost mai bine să creez o variantă nouă în care am șters codurile de nonrăspuns, păstrând-o și pe cea inițială. Să nu confundați ceea ce povestesc aici cu definirea nonrăspunsurilor din capitoul precedent. Definirea nonrăspunsurilor presupune că am păstrat codurile lor, numai că le dezactivăm din analizele pe care le rulăm. Eu vorbesc despre ștergerea fizică din bază.

O altă situație în care putem folosi **Recode into Different Variables** este atunci când vrem să prezentăm un tabel care conține încrucișarea dintre o variabilă măsurată metric, cum ar fi vârsta, și încrederea în oameni. Vârsta este înregistrată în ani împliniți : 18, 19, 20 etc. Încrederea în oameni este înregistrată folosind două variante de răspuns : „se poate avea încredere în cei mai mulți oameni” sau „e mai bine să fii atent în relațiile cu oamenii”. Dacă am realiza un tabel de contingență între cele două variabile, ar fi inutil, pentru că vârsta are foarte multe valori (în WVS 2012 pentru România, între 18 și 85 de ani). În această situație, alegem să recodificăm vârsta, adică să creăm o variabilă cu câteva categorii alese după un criteriu teoretic întemeiat stabilit de cercetător.



Am putea alege categoriile : 18-29, 30-39, 40-49, 50-59, 60+ . Observăm că pierdem informație. Aducem în aceeași categorie persoane de vârste diferite. De aceea categoriile nu se fac la întâmplare, ci motivat. Persoanele care sunt incluse în aceeași categorie trebuie să aibă trăsături comune, dar și diferite față de ale persoanelor din celelalte categorii. Dacă folosim **Recode into Same Variables**, pierdem informația detaliată și nu mai avem cum să o recuperăm. Dacă folosim **Recode into Different Variables**, creăm o variabilă nouă care conține categoriile respective și o păstrăm pe cea originală cu toată informația. Pe parcursul analizei, s-ar putea să ne trebuiască alte categorii de vârstă, să zicem 18-24, 25-29 etc. Putem să le obținem oricând, pentru că avem variabila inițială. Aici trebuie să facem trimitere către procesul de elaborare a întrebărilor din chestionar. Există cele patru niveluri de măsurare teoretice : nominal, ordinal, interval și raport. Sintetic, proprietățile acestora sunt prezentate în tabelul 5.1.

**Tabelul 5.1.** Niveluri de măsurare

|  | Nominal | Ordinal | Interval | Raport |
|--|---------|---------|----------|--------|
| <b>Categorii</b>                                       | Da      | Da      | Da       | Da     |
| <b>Categorii ordonate</b>                              |         | Da      | Da       | Da     |
| <b>Distanța dintre categoriile ordonate este egală</b> |         |         | Da       | Da     |
| <b>Număr</b>   |         |         |          | Da     |

Nivelul de măsurare cel mai slab din punct de vedere statistic este cel nominal. Sexul are două categorii : masculin și feminin. Suntem obișnuiți ca acestora să le fie atribuite codurile 1 și 2. Dar codurile acestea puteau fi foarte bine înlocuite cu 1001 și 23. Nu avea nici o importanță. Sexul feminin nu este pe locul 2, pentru că primește codul 2, după nici un criteriu. La fel, sexul masculin nu este pe locul 1, pentru că primește codul 1, tot după nici un criteriu. O discuție detaliată a acestor concepte este întâlnită în orice manual de metodologie cantitativă sau de statistică. Aș sublinia doar această idee : dacă puteți măsura o variabilă folosind un nivel de măsurare de raport, atunci faceți acest lucru. Dacă nu se poate utiliza un nivel de măsurare de raport, atunci căutați să folosiți unul de interval sau măcar ordinal. Dacă nici acest lucru nu este posibil, atunci folosiți un nivel nominal. Dintr-un număr putem face oricâte categorii și de orice fel dorim. Din categorii nu putem face numere. Am văzut deseori chestionare în care vârsta este înregistrată sub formă de categorii. Oricât de detaliate ar fi, tot se pierde informație. Aș măsura o variabilă care este metrică sub formă de categorii, doar dacă mă ajută să reduc numărul de non-răspunsuri. Dar, în chestionar, aș utiliza ambele variante. De exemplu, aș întreba care este venitul din luna trecută, lăsând posibilitatea să declare o sumă și apoi aș întreba și în ce categorie se încadrează. Evident, operatorul de teren, dacă

a aflat suma, o va încadra singur în categoria aferentă. Dar dacă nu a aflat suma, datorită refuzului, poate afla măcar categoria.

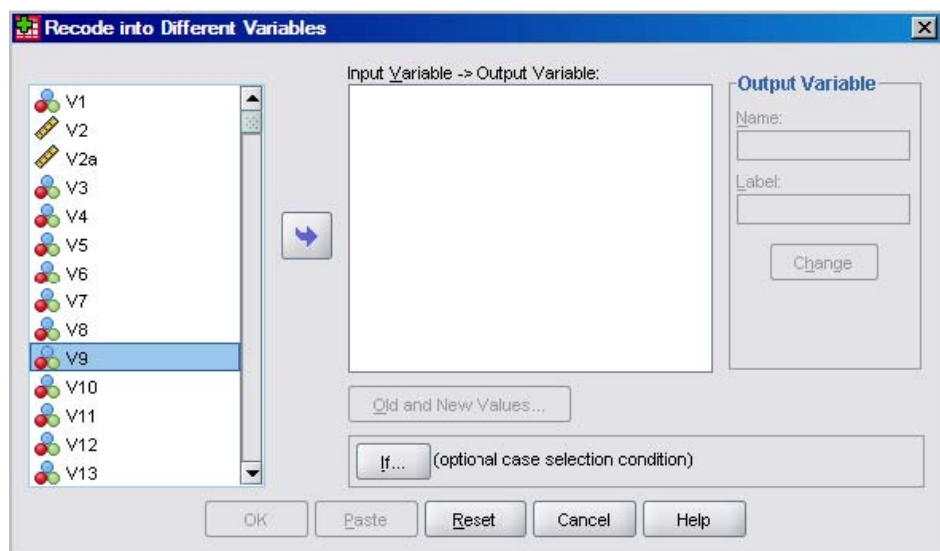
Variabilele dummy sunt un alt exemplu de utilizare a meniului **Recode into Different Variables**. O variabilă dummy ia valorile 1 sau 0. Valoarea 1 este atribuită caracteristicii care ne interesează, iar valoarea 0 celeilalte sau celorlalte. În regresia liniară muliplă nu pot folosi sexul codificat cu 1 și 2. Aleg cine primește 1, să zicem bărbații, iar 2 va fi transformat în 0. Dacă vreau să folosesc starea civilă ca predictor al fericirii și presupun că fenomenul explicat variază diferit pentru cei căsătoriți și pentru cei care au sau nu o relație, procedez astfel :

- presupun că cei căsătoriți sunt cei mai fericiți, așadar voi alege drept referință această categorie. Pentru ea nu mai creez un dummy ;
- creez un dummy în care 1 este atribuit celor care au o relație, dar nu sunt căsătoriți, iar 0 le este atribuit celor care nu au o relație, dar și celor căsătoriți ;
- creez un al doilea dummy în care 1 este atribuit celor care nu au o relație, fie sunt divorțați, fie sunt văduvi, iar 0 este atribuit celor care nu au o relație și celor care sunt căsătoriți.

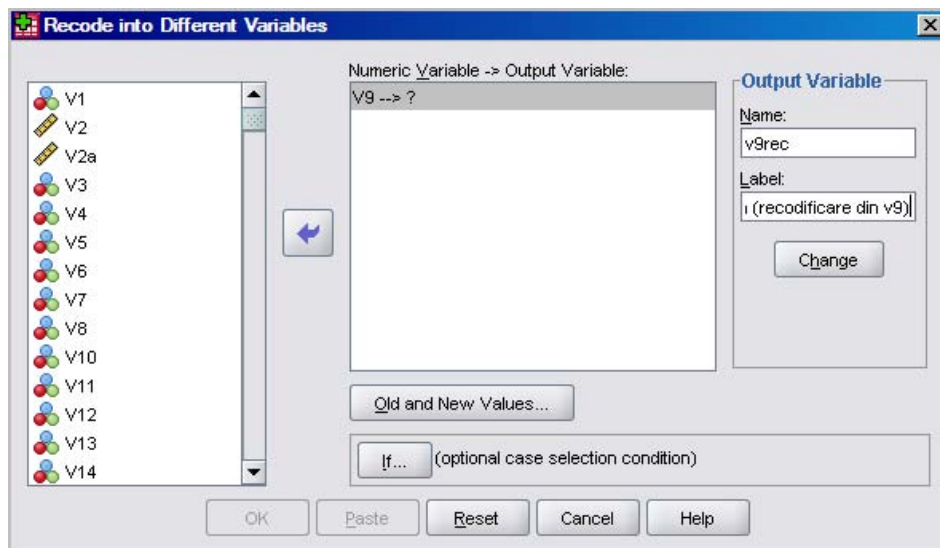
Voi prezenta meniul, folosind o altă situație care poate fi întâlnită în procesul de analiză : inversarea scalei de răspuns. În analiza elaborată de Voicu, Rusu și Comșa (2013), religiozitatea este măsurată ca gradul de importanță acordată religiei. În chestionarul folosit, întrebarea este „Vă rugăm să ne spuneți cât de importante sunt următoarele lucruri în viața dumneavoastră : ...religia ? ” și are patru variante de răspuns, de la „foarte importantă” (codul 1) la „deloc importantă” (codul 4). Pentru ca rezultatul analizei de regresie să fie mai ușor de citit, autorii au inversat scala astfel încât codul mare (4) să fie atribuit etichetei pozitive („foarte importantă”), iar codul mic (1) să fie atribuit etichetei negative („deloc importantă”). Măsura solidarității este orientată similar : un scor mare indică solidaritate ridicată. Astfel, dacă atunci când crește religiozitatea crește și solidaritatea, coeficienții de regresie vor avea semnul plus, iar interpretarea va fi intuitivă. Accesând meniul **Transform > Recode into Different Variables**, se deschide fereastra din figura 5.1a. Structura ferestrei ne este parțial familiară, pentru că seamănă cu cea de la **Recode into Same Variables**. Butonul **If** este folosit dacă dorim să punem o condiție care să fie activă atunci când creăm variabila nouă. În secțiunea **Output Variable**, care inițial este inactivă, introducem un nume pentru variabila pe care o creăm (**Name**) și o etichetă care explică numele (**Label**). Completarea informației la **Name** este obligatorie. La **Label** este opțională, dar recomandată. Altfel ar trebui să mergem în **Variable View** la coloana **Label** sau în sintaxă și să folosim comanda **VARIABLE LABELS**. Figura 5.1b prezintă cum se modifică interfața.

Figura 5.1. Meniul Transform &gt; Recode into Different Variables

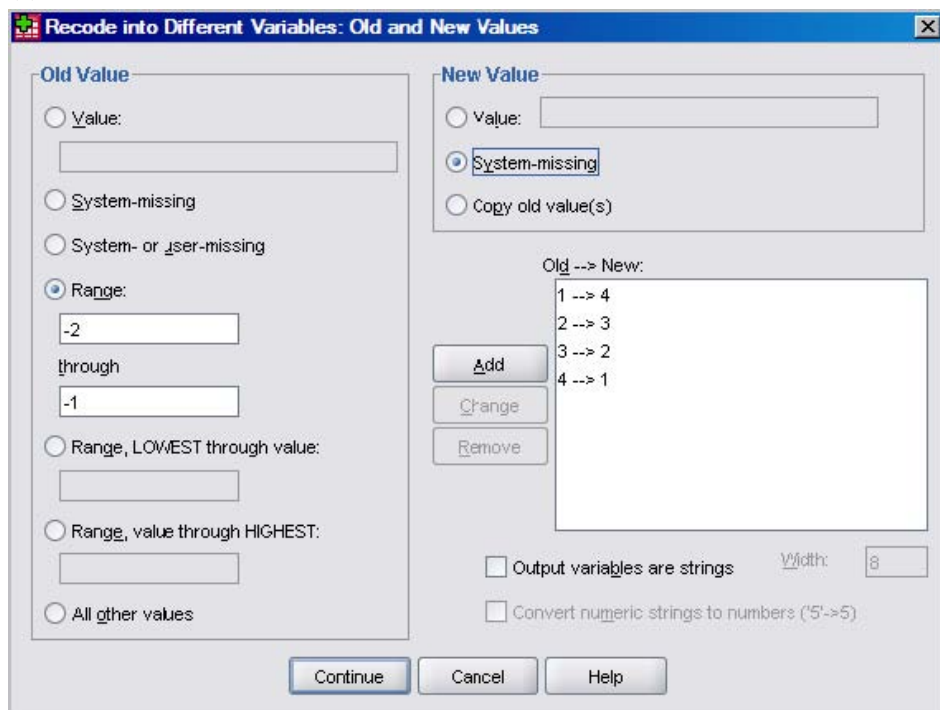
(a)



(b)



(c)



După ce introducem variabila V9 din lista de variabile din stânga în secțiunea **Numeric Variable --> Output Variable**, se activează **Name** și **Label**. Numele trebuie să respecte regulile programului. Eu prefer să adaug în coada numelui variabilei inițială expresia „rec” de la „recodificată”. Opțiunea dumneavoastră poate fi alta. La **Label** prefer să pun în etichetă informația „recodificare din variabila inițială”. Astfel, am o evidență clară a variabilelor pe care le-am creat. Odată completate aceste informații, trebuie să apăsăm butonul **Change**. Făcând acest lucru, dispare

semnul de întrebare (**V9 --> ?**) și apare :

Numeric Variable -> Output Variable:  
V9 --> v9rec

Următorul pas presupune să modificăm codurile conform nevoilor de analiză. Apăsăm butonul **Old and New Values** și se deschide fereastra din figura 5.1c. Fereastra are trei secțiuni: **Old Value**, **New Value** și cea care ne arată ce transformări facem. Mai întâi, trebuie să introducăm în fereastra **Old Value** codurile variabilei inițiale pe care le dorim transformate într-un fel sau altul. Aici dorim să inversăm scala: 1 devine 4, 2 devine 3, 3 devine 2, 4 devine 1. Se impune transformarea codurilor unul câte unul. O să lucrăm cu **Value**. După fiecare transformare, apăsăm butonul **Add**. Variabila V9 are și două coduri de nonrăspuns, -2 și -1. Pentru că nu vrem să le păstrăm în variabila nouă, le vom defini **System-missing**. Introducem la **Range** -2 la -1 și bifăm **System-missing**.

Variabila independentă pentru autoevaluarea stării de sănătate, care are numele V11 în baza de date, are distribuția din tabelul 5.2. O persoană nu a oferit un răspuns valid. Codurile sunt etichetate invers decât îmi doresc: codul mic (1) este asociat etichetei pozitive, iar codul mare (4) este asociat etichetei negative. Pentru că dorim să interpretăm efectul pozitiv al autoevaluării sănătății asupra satisfacției cu viața, așteptându-ne la o relație pozitivă (semn + la coeficientul de regresie), recodificăm variabila astfel încât codurile să fie în acord cu etichetele. De asemenea, o să ștergem din bază nonrăspunsul respectiv.

**Tabelul 5.2.** Tabel de frecvență pentru autoevaluarea stării de sănătate

| <b>V11 State of health (subjective)</b> |              |           |         |               |                    |
|---|--------------|-----------|---------|---------------|--------------------|
|   |              | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                                   | 1 Very good  | 232       | 15.4    | 15.4          | 15.4               |
|   | 2 Good       | 747       | 49.7    | 49.7          | 65.2               |
|   | 3 Fair       | 390       | 25.9    | 25.9          | 91.1               |
|   | 4 Poor       | 134       | 8.9     | 8.9           | 100.0              |
|   | Total        | 1502      | 99.9    | 100.0         |                    |
| Missing                                 | -2 No answer | 1         | .1      |               |                    |
| Total                                   |              | 1503      | 100.0   |               |                    |

Așadar, mergem în meniul **Transform > Recode into Different Variables**. Dacă nu am închis baza de date între timp, o să observăm comenzile de la recodificarea anterioară. Pentru a nu ne complica, apăsăm butonul **Reset**, acesta ștergând orice informație care era prezentă în meniul respectiv. Trecem V11 în dreapta. Îi dăm un nume: v11rec. Etichetăm numele: autoevaluarea stării de sănătate (recodificare din V11). Apăsăm **Change**. Apăsăm **Old and New Values**. Transformăm pe rând fiecare valoare folosind **Value** de la **Old Value**: 1 --> 4, 2 --> 3, 3 --> 2, 4 --> 1, -2 --> **System-missing**. Apăsăm butonul **Add** după fiecare modificare. **Continue** și **OK**.

După recodificare trebuie să verificăm dacă am lucrat corect. În cazul creării unei noi variabile prin recodificare, vom realiza un tabel de contingență (**Crosstab**) dintre variabila inițială (V9 sau V11) și variabila nou-creată (v9rec sau v11rec). Acest tabel se realizează din meniul **Analyze > Descriptive Statistics > Crosstabs**. Pe rând (**Row**) introducem variabila cu mai multe categorii. Pe coloană (**Column**) introducem variabila creată (tabelul 5.3).

În primul rând, observăm că nu avem etichete pentru codurile variabilei nou-create: 1, 2, 3 și 4. Deci trebuie să le introducem în coloana **Values** din **Variable View** sau folosind sintaxa de mai jos. Trebuie reținut că acesta este, de cele mai multe ori, primul pas după recodificare.

#### **VALUE LABELS** v9rec

- 1 not at all important
- 2 not very important
- 3 rather important
- 4 very important

După ce am rulat această sintaxă, realizăm din nou tabelul. Rezultatul este vizibil în tabelul 5.3b. Acum este mai ușor de citit. Al doilea lucru pe care îl observăm este că nu mai apar coduri de nonrăspuns. Dacă lucrăm cu **user-missing** sau **system missing**, comanda **Crosstabs** le va ignora. Pe noi ne interesează, în acest tabel, să vedem dacă etichetelor le corespund oamenii potriviți. Aceștia sunt distribuiți pe diagonală, deci am lucrat corect. Atenție însă: dacă am etichetat greșit, programul nu ne va avertiza. Să fim atenți la fiecare etapă de lucru.

**Tabelul 5.3.** Tabel de contingență pentru verificarea corectitudinii recodificării

(a)

| <b>V9 Important in life: Religion * v9rec importanta religiei in viata (recodificare din v9) Crosstabulation</b> |                        |  |     |     |     |       |
|--|------------------------|--|-----|-----|-----|-------|
| Count  |                        |  |     |     |     |       |
|  |                        | v9rec importanta religiei în viață (recodificare din v9) |     |     |     | Total |
|  |                        | 1  | 2   | 3   | 4   |       |
| V9 Important in life: Religion   | 1 Very important       | 0  | 0   | 0   | 758 | 758   |
|  | 2 Rather important     | 0  | 0   | 500 | 0   | 500   |
|  | 3 Not very important   | 0  | 192 | 0   | 0   | 192   |
|  | 4 Not at all important | 48   | 0   | 0   | 0   | 48    |
| Total  |                        | 48   | 192 | 500 | 758 | 1498  |

(b)

| <b>V9 Important in life: Religion * v9rec importanta religiei în viață (recodificare din v9) Crosstabulation</b> |                        |  |                      |                    |                  |       |
|--|------------------------|--|----------------------|--------------------|------------------|-------|
| Count  |                        |  |                      |                    |                  |       |
|  |                        | v9rec importanta religiei in viata (recodificare din v9) |                      |                    |                  | Total |
|  |                        | 1 not at all important                                   | 2 not very important | 3 rather important | 4 very important |       |
| V9 Important in life: Religion   | 1 Very important       | 0  | 0                    | 0                  | 758              | 758   |
|  | 2 Rather important     | 0  | 0                    | 500                | 0                | 500   |
|  | 3 Not very important   | 0  | 192                  | 0                  | 0                | 192   |
|  | 4 Not at all important | 48   | 0                    | 0                  | 0                | 48    |
| Total  |                        | 48   | 192                  | 500                | 758              | 1498  |

## 5.2. Crearea unei alte variabile utilizând meniul Compute

Am discutat deja o situație în care folosim **Transform > Compute**. Pentru că mi se pare important, am să insist prezentând, pentru început, cum realizăm o variabilă de ponderare (**weight**).

Să presupunem că ponderăm în funcție de mediul de rezidență, vârstă și sex. Mai întâi trebuie să stabilim care sunt categoriile pentru fiecare criteriu. Categoriile se aleg și în funcție de cum este disponibilă informația pentru acestea. Am ales categoriile urban și rural pentru mediul de rezidență și categoriile 18-34, 35-49, 50-64, 65+ pentru vârstă. Pentru sex avem doar două categorii: bărbat sau femeie. Căutăm la Institutul Național de Statistică informații pentru tabelul :

| Vârstă | Bărbați în Urban | Femei în Urban | Bărbați în Rural | Femei în Rural | Total |
|--------|------------------|----------------|------------------|----------------|-------|
| 18-34  | ...              | ...            | ...              | ...            | ...   |
| 35-49  | ...              | ...            | ...              | ...            | ...   |
| 50-64  | ...              | ...            | ...              | ...            | ...   |
| 65+    | ...              | ...            | ...              | ...            | ...   |
| Total  | ...              | ...            | ...              | ...            | ...   |

Mai concret, informațiile pe care trebuie să le punem în fiecare celulă sunt numerele de persoane care se încadrează simultan în toate cele trei categorii desemnate de rândurile și coloanele tabelului. De exemplu, bărbații care locuiesc în urban și au vârsta cuprinsă între 18-34 de ani ar putea fi în număr de 1.600.000. Realizăm același tabel și pentru eșantion. Evident, numerele din fiecare celulă vor fi mult mai mici, dată fiind mărimea eșantionului. De exemplu, în eșantion ar putea fi incluși 106 bărbați care locuiesc în urban și au vârsta cuprinsă în intervalul 18-34 de ani

Calculăm proporția fiecărei celule din totalul populației, respectiv a eșantionului. Vor rezulta două noi tabele care conțin aceste proporții. Apoi vom împărți proporțiile din populație la proporțiile din eșantion :

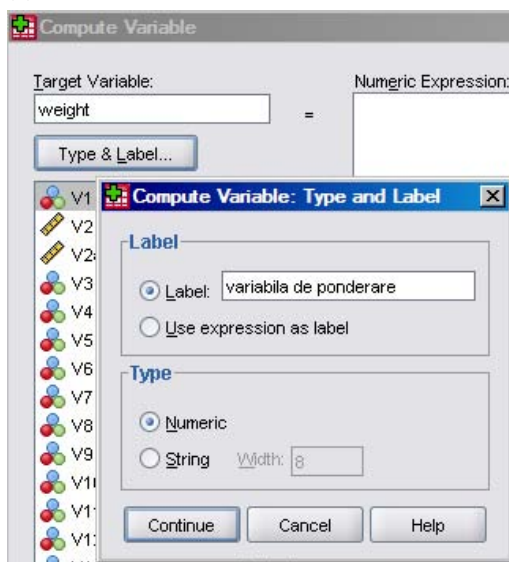
| Vârstă | Bărbați în Urban | Femei în Urban | Bărbați în Rural | Femei în Rural | Total |
|--------|------------------|----------------|------------------|----------------|-------|
| 18-34  | 9,16/7.07=1,2942 | ...            |                  |                |       |
| 35-49  | ...              |                |                  |                |       |
| 50-64  |                  |                |                  |                |       |
| 65+    |                  |                |                  |                |       |
| Total  |                  |                |                  |                |       |

Rezultatul final reprezintă valorile pe care le va lua ponderea pentru fiecare dintre aceste categorii compuse. Această nouă variabilă trebuie introdusă în SPSS. Realizăm acest lucru cu meniul **Transform > Compute**. În secțiunea **Target**

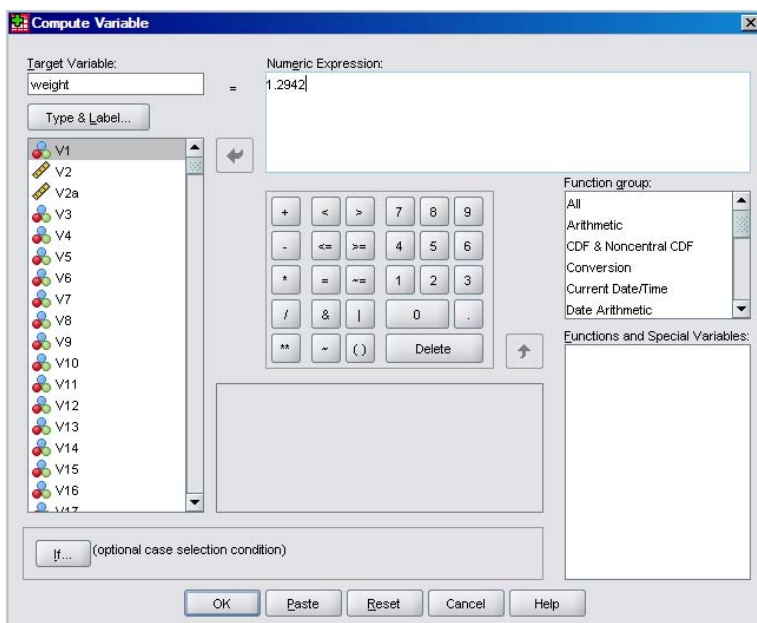
**Variable**, dăm un nume variabilei pe care o realizăm, căreia îi atribuim și o etichetă în secțiunea **Label** din fereastra care se deschide apăsând butonul **Type & Label** (figura 5.2a).

**Figura 5.2.** Transform > Compute : crearea unei variabile de ponderare

(a)

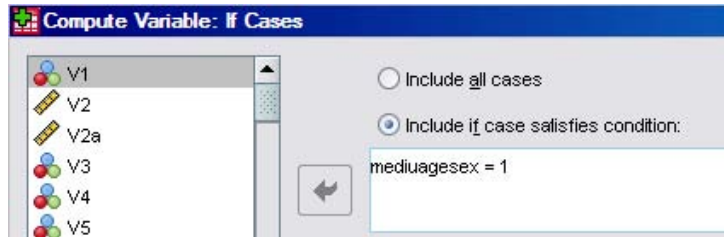


(b)





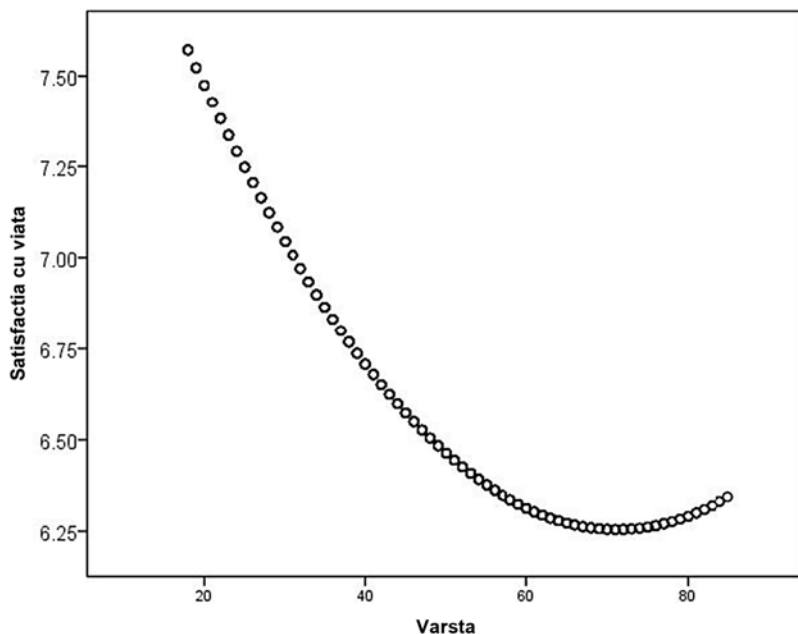
(c)



Observăm semnul „=” în dreapta câmpului **Target Variable**. În dreapta acestui semn, în câmpul **Numeric Expression**, introducem formula prin care realizăm noua variabilă. În cazul de față, nu avem o formulă: doar imputăm valoarea 1.2942 pe care o introducem fie din tastatură, fie folosind butoanele din centrul ferestrei (figura 5.2b). Dacă apăsăm **OK** acum, variabila de ponderare (**weight**) va avea valoarea 1.2942 pentru toate persoanele din eșantion. Însă această pondere este doar pentru categoria bărbaților care locuiesc în urban și au vârsta în intervalul 18-34 de ani. De aceea trebuie să folosim și butonul **If...** din colțul stânga jos al ferestrei. Apăsând acest buton se deschide fereastra din figura 5.2c.

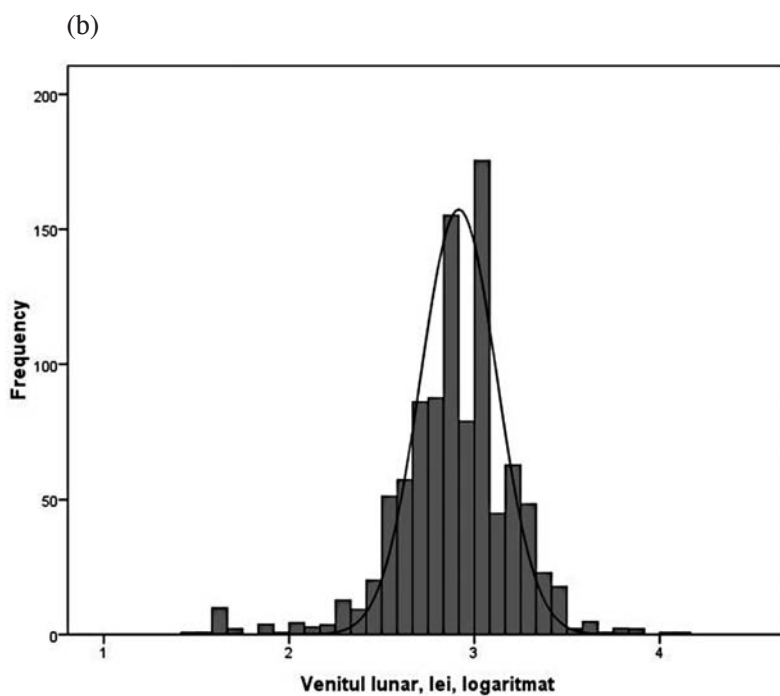
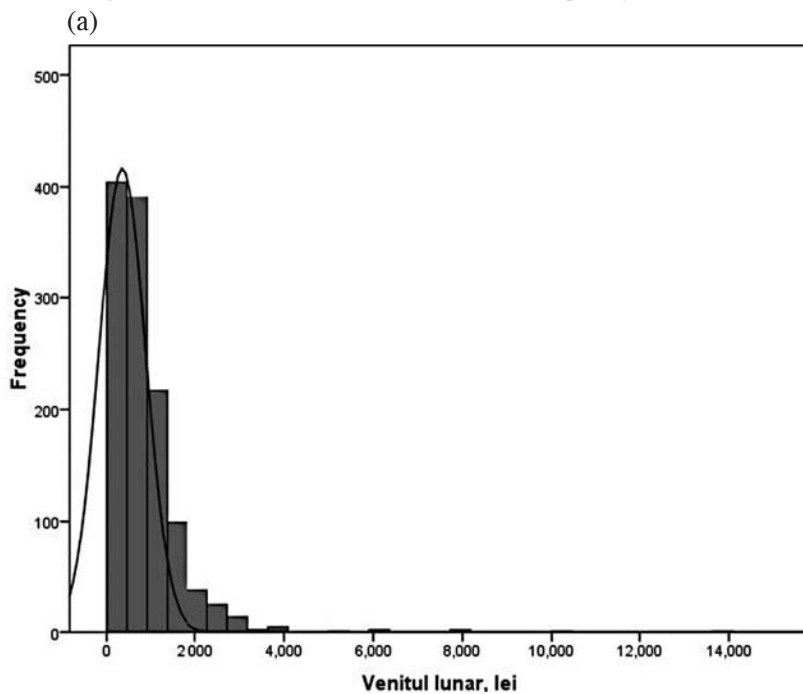
Inițial este bifată opțiunea **Include all cases**. Pentru că vrem să punem o condiție, vom bifa **Include if case satisfies condition** și vom introduce condiția în câmpul activat. În exemplul nostru, am presupus că în baza de date există deja o variabilă care reflectă apartenența simultană la cele trei categorii. Aceasta are numele **mediuagesex** și conține 16 categorii. Codul 1 reprezintă categoria bărbați care locuiesc în mediul urban și au vârsta între 18 și 34 de ani. Pentru că ponderea 1.2942 este ponderea pentru această categorie, atunci vom introduce aici condiția **mediuagesex = 1**. Astfel SPSS va atribui ponderea 1.2942 doar categoriei 1 de la variabila **mediuagesex**. Repetăm procedura pentru toate celelalte categorii.

Un alt exemplu. Relația dintre satisfacția cu viața și vârstă nu este liniară (Lelkes, 2008). Adică satisfacția nu crește/decrește, constant, odată cu înaintarea în vârstă. Mai degrabă, cele două au o relație nonlineară asemănătoare cu cea reprezentată în figura 5.3. Cel mai înalt nivel al satisfacției cu viața este trăit în tinerețe, când grijile materiale și sociale nu sunt atât de multe, părinții încă îi întrețin pe copii etc. Urmează momente cum ar fi cel al intrării pe piața muncii, al formării propriei familii, al accentuării independenței financiare etc. Copiii pleacă de acasă, grijile cu privire la siguranța locului de muncă se accentuează etc. Vine vârsta pensionării, grijile legate de profesie se reduc, dar apar probleme de sănătate asociate vârstei, moartea partenerului de viață etc. Pe de altă parte, oamenii își pot urmări interesele personale mai mult decât înainte, cel puțin prin prisma timpului liber de care dispun. Toate acestea sunt explicații plauzibile pentru acest tip de relație dintre vârstă și satisfacția vieții.

**Figura 5.3.** Relație nonliniară dintre vârstă și satisfacția cu viața

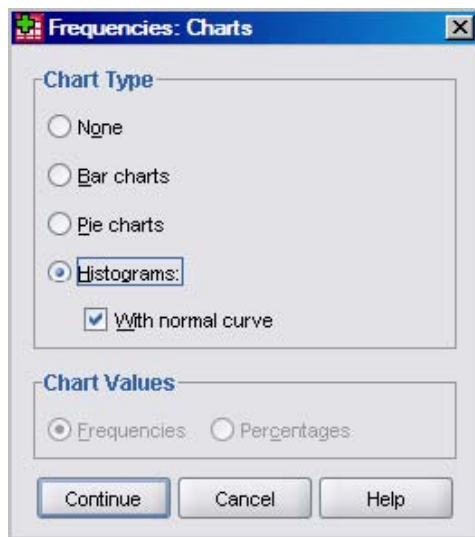
Dacă vrem să aplicăm o modelare liniară, trebuie să includem și acel punct de inflexiune în analiză. Acest lucru se face aici prin introducerea variabilei vârstă ridicată la pătrat, alături de variabila vârstă inițială. Variabila vârstă la pătrat va fi creată folosind meniul **Compute**. La **Numeric Expression** notăm formula :  $\text{varsta2} = \text{V242} * \text{V242}$ . V242 este variabila care conține vârsta respondentului din WVS 2012.

Astfel de transformări sunt frecvente în analizele multivariate. Una care folosește funcțiile implementate în SPSS presupune calcularea unui logaritm. Această transformare este frecvent întâlnită pentru variabila venit care nu are o distribuție normală, ci, de regulă, alungită la dreapta (figura 5.4a). Majoritatea românilor au venituri mici, dar există și români care au venituri ceva mai mari. Unii dintre aceștia pot să se îndepărteze destul de mult de majoritate. În analizele statistice, aceștia sunt considerați cazuri extreme (**outlieri**). Trebuie văzut în ce măsură afectează rezultatele analizelor statistice pe care le rulăm. Putem transforma variabila folosind una dintre funcțiile de logaritmare. La **Numeric Expression** aducem din secțiunea **Functions and Special Variables**, dând dublu click pe ea, funcția **LG10()**. Trebuie doar să introducem între paranteze, în locul semnului de întrebare, variabila din baza de date care conține informații despre venit : **cs237a** în WVS 2012. Funcția devine **LG10(cs237a)**. Apăsăm **OK**. Distribuția variabilei logaritmă aproximează mai bine forma așteptată (figura 5.4b). Problema acestor transformări este creșterea gradului de dificultate a interpretării coeficienților de regresie atunci când, în locul unității de măsură a variabilei inițiale, folosim logaritmi sau rezultatele altor funcții matematice.

**Figura 5.4.** Distribuția venitului înainte și după logaritmare

Histograma din figura 5.4 a fost creată din meniul **Analyze / Descriptive Statistics / Frequencies**. În fereastra care s-a deschis, apăsăm pe butonul **Charts** (figura 5.5). Inițial este selectat **None**, dar noi suntem interesați de histogramă, de aceea facem selecția corespunzătoare : **Histogram > With normal curve**.

**Figura 5.5.** Realizarea graficelor din meniul Frequencies > Charts



### 5.3. Exerciții

Pentru aceste exerciții utilizăm baza de date și/sau chestionarul World Values Survey 2012 rezultat(ă)e în urma aplicării chestionarului în România. Baza de date poate fi descărcată de pe pagina de internet a *Grupului Românesc pentru Studiul Valorilor Sociale* (<http://www.romanianvalues.ro>).

1. Căutați pe siteul [www.romanianvalues.ro](http://www.romanianvalues.ro) newsletterul nr. 4 din 2013-2014 cu tema „Satisfacția cu viața”. Citiți acest text și realizați o listă cu variabilele utilizate în analiză.
2. Găsiți variabila evaluarea stării de sănătate. Creați o variabilă dummy pornind de la aceasta. Căror coduri le atribuiți valoarea 1 și căror coduri le atribuiți valoarea 0? Argumentați decizia.
3. Găsiți variabila stare civilă. Creați o variabilă dummy pornind de la aceasta, astfel încât să reflecte categoriile „persoana are o relație” / „persoana nu are o relație”.
4. Care sunt variabilele dummy pe care le puteți crea din punct de vedere teoretic pornind de la variabila stare civilă? Este fezabil să le creați pe toate? Argumentați răspunsul.

5. Creați o variabilă care să conțină următoarele categorii de vârstă: 18-24, 25-29, 30-34, 35-39, 40-44, 45-49, 50-54, 55-59, 60-64, 65+.
6. Căutați în baza de date setul de variabile v102-v107 care se referă la încrederea în diferite categorii de persoane. Recodificați, creând variabile noi, toate aceste variabile astfel încât din patru variante de răspuns să rămâneți doar cu două.
7. Citiți lista de variabile v96-v101 din chestionar. Creați o nouă variabilă care să reprezinte suma tuturor acestor variabile. Ce măsoară această variabilă? Cum interpretați scorul 60?
8. Căutați pe siteul [www.romanianvalues.ro](http://www.romanianvalues.ro) newsletterul nr. 2 din 2013-2014 cu tema „Încrederea în instituții”. Citiți secțiunea „Cine are încredere în instituțiile politice?”. Identificați în baza de date variabilele sex, vârstă, nivel de educație, mediu de rezidență, autopoziționare în clasa socială și mândria de a fi român. Recodificați aceste variabile astfel încât să corespundă modului în care sunt utilizate în această lucrare.



## 6. O primă privire asupra datelor

Prin cercetările noastre, căutăm să descriem și/sau să explicăm un anumit fenomen social. Câți români au emigrat în anul 2013? Dintre aceștia, câți au studii superioare? Câți români suferă de o boală cronică? Dintre aceștia, câți au vârsta cuprinsă între 18 și 30 de ani? Câți români consumă pufuleți? Dintre aceștia, câți cumpără pufuleții din hipermarket și câți din magazinul din proximitatea locuinței? În primul rând, descriem situația, dar, de regulă, vrem să și explicăm de ce situația arată în felul acesta.

Testăm prezența unui efect (dacă..., atunci...), dar și intensitatea cu care variază o variabilă dependentă în funcție de variația variabilei independente (cu cât..., cu atât...). Putem compara nivelul de satisfacție cu viața al românilor care au emigrat cu cel al românilor care au decis să nu facă acest lucru. Dacă există diferențe, emigrarea este factorul care produce diferența sau pot fi identificați și alți factori? Nivelurile de satisfacție cu viața sunt similare în cazul tuturor celor ce au emigrat sau variază în funcție de caracteristicile țării de destinație? Care este factorul care crește cel mai mult satisfacția cu viața? Sunt mai satisfăcuți cu viața cei care au emigrat când erau mai tineri (sub 25 de ani) sau cei care au emigrat la o vârstă mai înaintată (peste 25 de ani)?




Primul pas în acest demers este să ne familiarizăm cu datele. Să descriem modul în care gândește și se comportă majoritatea. Primul pas este analiza statistică univariată. Avem o listă de variabile aleasă conform obiectivelor de cercetare și, pentru fiecare dintre acestea, inspectăm distribuțiile și diferiți indicatori statistici care pot fi calculați pentru ele. Citim datele într-o manieră descriptivă. Dacă ne interesează să identificăm motivele pentru care unii români sunt mai fericiți decât alții sau motivele pentru care unii români își autoevaluează sănătatea ca fiind mai bună decât a altora, atunci începem prin a ne uita la distribuția fericirii sau stării de sănătate a românilor la momentele alese pentru perspectiva cercetării. Câți români sunt fericiți și câți nefericiți? După care, trecem la analizele statistice bivariate. Începem să punem în relație variabilele din lista noastră, două câte două. Care sunt categoriile care cuprind cei mai mulți români fericiți: locuitorii orașelor mici sau ai orașelor mari, tinerii sau adulții, cei căsătoriți sau persoanele care nu au o relație de cuplu, cei care au absolvit facultatea sau cei care au absolvit doar liceul, cei din cuartila unu, doi sau trei de venit ș.a.m.d.? Câți români apreciază că starea lor de sănătate este bună și câți o apreciază ca fiind proastă? Fenomenele sociale sunt complexe, de aceea analizele uni- sau bi-variate sunt insuficiente pentru a înțelege adecvat variația acestora. Orice analist dorește

să ajungă la analizele statistice multivariate. Cine sunt cei care își evaluează sănătatea ca fiind mai bună: vegetarienii sau omnivorii, cei care merg la medic pentru controale preventive, cei mai educați, cei care fac sport ș.a.m.d. ?

Descrierea datelor se realizează prin calcularea unor indicatori statistici și, vizual, prin inspectarea unor grafice. Calculăm indicatori ai tendinței centrale, media (**mean**) și mediana (**median**), dar și indicatori ai variației, abaterea standard (**standard deviation**) sau coeficientul de variație. Realizăm grafice bară (**bar chart**), histogramă (**histogram**) sau nor de puncte (**scatterplot**).

## 6.1. Cum gândește majoritatea și cât de omogene sunt grupurile comparate

Indicatorii sintetici, cum sunt media sau mediana, oferă rapid, printr-un singur număr, o imagine de ansamblu asupra situației majorității din populația de referință. Alții, cum este abaterea standard, ne arată cât de omogene sunt, după aceeași caracteristică, diferite grupuri. Media și mediana sunt indicatori ai tendinței centrale. Abaterea standard este un indicator al variației.

Acești indicatori pot fi calculați doar atunci când variabilele au anumite proprietăți. Aceste proprietăți sunt grupate sub numele de niveluri de măsurare (tabelul 5.1). Mediana este valoarea care împarte setul de date ordonate în două părți egale. Poate fi calculată dacă variabila are cel puțin nivelul de măsurare ordinal sau, în limbajul cercetătorilor, este variabilă ordinală. Media poate fi calculată doar pentru variabile metrice, interval sau raport. Pentru variabilele nominale, vom inspecta distribuția de frecvențe: categoria cu cele mai multe unități va fi tendința centrală. Dacă ne reamintim coloanele din **Variable View**, mai exact coloana **Measure**, remarcăm că SPSS distinge între variabilele nominale ( **Nominal**), ordinale ( **Ordinal**) și metrice ( **Scale**). În cercetarea socială, atunci când aplicăm un chestionar, este destul de greu să măsurăm prin procedeele uzuale, la nivel de raport. În cel mai fericit caz, am reușit să elaborăm variabile ordinale sau de interval. De aceea, în practică, pentru interval și raport sunt folosite aproximativ aceleași analize statistice. O discuție care clarifică multe dintre aceste aspecte este oferită de Agresti și Finlay (2008).

În științele sociale, folosim frecvent media aritmetică pentru a reprezenta tendința centrală. Este larg cunoscută, majoritatea știind să o interpreteze. Spre deosebire de mediană, utilizează informația numerică din variabilă, nu doar ordinea scorurilor (Agresti și Finlay, 2008). Totuși, mediana este frecvent consultată de analist: este, cel puțin, un mecanism de verificare a mediei sau chiar înlocuitor al acesteia, atunci când datele conțin cazuri extreme (**outliers**). Cazurile extreme sunt persoane care au valori mult mai mari sau mult mai mici decât majoritatea la variabila respectivă. O persoană care are un salariu lunar de 25.000



de lei, în condițiile în care următorul salariu, în ordine descendentă, este de 8.000 de lei, este un caz extrem. Aceasta nu este o situație ireală. Ea are însă un impact negativ asupra calculelor statistice. Prezența printre valorile variabilei chiar și a unui singur caz extrem, indiferent că se află în partea stângă (valoare foarte mică) sau în partea dreaptă a scalei (valoare foarte mare), va afecta serios media, micșorându-i sau crescându-i foarte mult valoarea. Calculând salariul mediu folosind și valoarea 25.000 lei va distorsiona media: salariul mediu va lua o valoare care nu reflectă situația majorității. Rotariu, Bădescu și colaboratorii (2006), prezentând detaliat proprietățile mediei și medianei, atrag atenția că media nu este valoarea mijlocie a seriei. Media se va încadra în intervalul valorilor variabilei pentru care este calculată, fiind exprimată în aceeași unitate de măsură cu aceasta. Dacă variabila este „salariu exprimat în lei”, atunci media va fi exprimată în lei. De Vaus (2002), la rândul său, subliniază un alt neajuns al mediei, care trebuie avut în considerare în momentul interpretării valorii calculate de program: aceeași medie poate fi obținută din distribuții diferite. Agresti și Finlay (2008) demonstrează cum media este deplasată în direcția cozii mai lungi, atunci când distribuția este alungită la stânga sau la dreapta. Când grupurile pentru care este calculată sunt omogene, adică persoanele seamănă între ele, media va fi un indicator bun al tendinței centrale, dar mai puțin bun atunci când grupurile sunt eterogene. Acesta este unul dintre motivele pentru care calculăm și indicatori ai variației sau dispersiei, împreună cu indicatorii tendinței centrale.

Indicatorii variației sau dispersiei arată gradul de împrăștiere sau omogenitate/eterogenitate a grupurilor investigate după o variabilă anume. Înainte de a calcula un indicator al variației, trebuie să stabilim ce nivel de măsurare are variabila respectivă. Cel mai utilizat indicator este abaterea standard, care, pentru că folosește media în formula de calcul, poate fi calculat doar pentru variabile metrice. Putem compara abaterile standard calculate pentru aceeași variabilă în cazul a două grupuri. Grupul care arată cea mai mare abatere standard va fi mai eterogen. Dar această comparație nu ne va spune prea multe despre cât de omogen sau eterogen este fiecare grup. Agresti și Finlay (2008) prezintă o regulă empirică aplicabilă distribuțiilor aproximativ normale, pe care o putem utiliza pentru a interpreta abaterea standard și în termenii mărimii valorii acesteia: (1) aproximativ 68% dintre cazuri se află în intervalul [medie – abatere standard, medie + abatere standard], (2) aproximativ 95% dintre cazuri se află în intervalul [medie – 2 x abatere standard, medie + 2 x abatere standard] și (3) aproape toate cazurile se află în intervalul [medie – 3 x abatere standard, medie + 3 x abatere standard]. Abaterea standard are câteva neajunsuri care pot fi corectate prin utilizarea altui indicator al variației, coeficientul de variație. Coeficientul de variație este egal cu raportul dintre abaterea standard și media variabilei. Acesta este util atunci când vrem să comparăm anumite grupuri (1) folosind o variabilă care are unități de măsură diferite și/sau (2) nivelul general al valorilor variabilei este diferit în grupurile respective. Rotariu, Bădescu și colaboratorii (2006) oferă o explicație detaliată în acest sens: nu poți compara salariile din România, exprimate în lei, cu cele din Germania, exprimate în euro, la fel cum nu poți compara masa corporală

a unor albine cu cea a unor elefanți. Acești autori atrag atenția la utilizările fără logică teoretică ale coeficientului de variație: poate fi calculat doar pentru nivelul de măsurare de raport, pentru că valorile au originea zero. De asemenea, nu trebuie utilizat pentru a compara grupurile folosind variabile care au conținut diferit.

Sunt situații în care dorim să știm ce procent din observații se află sub sau deasupra unei valori. Acest gen de informație ne este oferit, de exemplu, de mediană: 50% dintre observații se află sub această valoare și 50% peste această valoare. Pentru informații mai detaliate utilizăm percentilele, întâlnite în cărțile de statistică sub denumirea de măsuri ale poziționării (Agresti și Franklin, 2013). Percentilele sunt de mai multe feluri. Cuartilele sunt foarte utilizate. Există trei cuartile, cuartila 2 fiind chiar mediana. Sub prima cuartilă se află 25% dintre cazuri, iar deasupra celei de-a treia cuartile se află tot 25% dintre cazuri. Cel mai simplu este să vă reprezentați o linie împărțită în patru segmente, fiecare segment reprezentând 25% din date. Asociată cuartilelor este abaterea intercuartilă, care ne arată distanța dintre cuartilele trei și unu. Din acest motiv, abaterea intercuartilă nu este influențată de cazurile extreme, fiind utilizată pentru detectarea acestora: dacă o observație se află dincolo de  $1.5 \times \text{AIQ}$ , adică sub prima cuartilă sau peste a treia cuartilă, atunci s-ar putea să fie un caz extrem. Graficul **box-plot** ne ajută să vizualizăm acest gen de informații.

SPSS oferă mai multe posibilități prin care putem calcula indicatorii tendinței centrale, variației și poziționării.

Pentru variabilele nominale, utilizăm distribuția de frecvențe pe care o obținem din meniul **Analyze > Descriptive statistics > Frequencies**. La întrebarea „În general vorbind, ați spune că se poate avea încredere în cei mai mulți oameni sau că e mai bine să fii atent în relațiile cu oamenii?” adresată în WVS 2012 și românilor, distribuția răspunsurilor este cea prezentată în tabelul 6.1. În primul rând, remarcăm cele 15 persoane care nu au oferit un răspuns valid (coloana **Frequency**) (tabelul 6.1a). Trebuie să instruim programul că -2 și -1 sunt coduri de nonrăspuns care trebuie dezactivate din analiză. Facem acest lucru fie în coloana **Missing** din **Variable View** (**Discrete missing values = -2, respectiv -1**), fie rulând sintaxa **MISSING VALUES V24 (-2, -1)**. Rezultatul este prezentat în tabelul 6.1b.

**Tabelul 6.1.** Tabel de frecvență: înainte și după definirea nonrăspunsurilor

(a)

| <b>V24 Most people can be trusted</b> |                              |           |         |               |                    |
|---------------------------------------|------------------------------|-----------|---------|---------------|--------------------|
|                                       |                              | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                                 | -2 No answer                 | 4         | .2      | .2            | .2                 |
|                                       | -1 Don't know                | 11        | .7      | .7            | 1.0                |
|                                       | 1 Most people can be trusted | 115       | 7.7     | 7.7           | 8.6                |
|                                       | 2 Need to be very careful    | 1373      | 91.4    | 91.4          | 100.0              |
|                                       | Total                        | 1503      | 100.0   | 100.0         |                    |

(b)

| <b>V24 Most people can be trusted</b> |                              |           |         |               |                    |
|---------------------------------------|------------------------------|-----------|---------|---------------|--------------------|
|                                       |                              | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                                 | 1 Most people can be trusted | 115       | 7.7     | 7.7           | 7.7                |
|                                       | 2 Need to be very careful    | 1373      | 91.4    | 92.3          | 100.0              |
|                                       | Total                        | 1489      | 99.0    | 100.0         |                    |
| Missing                               | -2 No answer                 | 4         | .2      |               |                    |
|                                       | -1 Don't know                | 11        | .7      |               |                    |
|                                       | Total                        | 14        | 1.0     |               |                    |
| Total                                 |                              | 1503      | 100.0   |               |                    |

Observăm că 92 % dintre români considerau, în 2012, că e mai bine să fii atent în relațiile cu oamenii.

În același an, majoritatea românilor considerau că principala problemă din lume este sărăcia : 53 % au ales această variantă de răspuns în defavoarea celorlalte (tabelul 6.2).

**Tabelul 6.2.** Tabel de frecvență : după definirea nonrăspunsurilor

| <b>V80 Most serious problem of the world</b> |   |           |         |               |                    |
|--|---|-----------|---------|---------------|--------------------|
|  |   | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid  | 1 People living in poverty and need       | 782       | 52.0    | 52.7          | 52.7               |
|  | 2 Discrimination against girls and women  | 105       | 7.0     | 7.1           | 59.8               |
|  | 3 Poor sanitation and infectious diseases | 205       | 13.7    | 13.9          | 73.7               |
|  | 4 Inadequate education                    | 260       | 17.3    | 17.5          | 91.2               |
|  | 5 Environmental pollution                 | 130       | 8.7     | 8.8           | 100.0              |
|  | Total                                     | 1483      | 98.7    | 100.0         |                    |
| Missing                                      | -2 No answer                              | 7         | .5      |               |                    |
|  | -1 Don't know                             | 12        | .8      |               |                    |
|  | Total                                     | 20        | 1.3     |               |                    |
| Total  |   | 1503      | 100.0   |               |                    |

Remarcați diferența dintre coloana **Percent** și **Valid Percent**. În prima sunt calculate procentele luând ca bază întregul eșantion, adică și pe cei care nu au

oferit un răspuns valid. În cea de-a doua sunt calculate procentele luând ca bază eșantionul valid, adică doar pe cei care au oferit un răspuns valid.

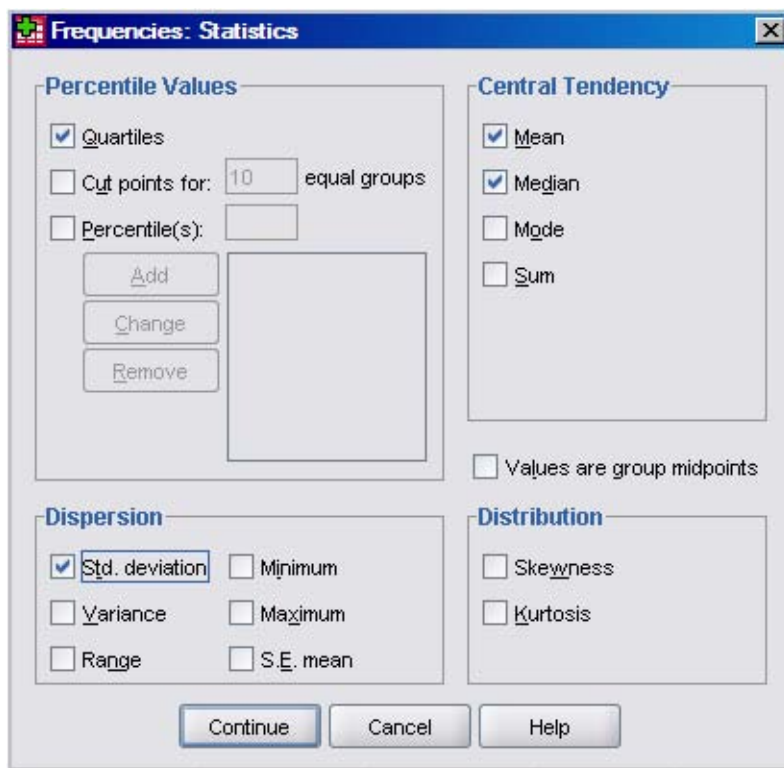
Pentru indicatorii tendinței centrale, variației și poziționării, calculabili pentru variabilele care au cel puțin nivelul de măsurare ordinal, putem utiliza meniurile **Analyze > Descriptive Statistics > Frequencies**, **Analyze > Descriptive Statistics > Descriptives** sau **Analyze > Descriptive Statistics > Explore**.

Meniul **Analyze > Descriptive Statistics > Frequencies** ne este deja familiar pentru că l-am folosit pentru a realiza tabelele de frecvență. Până acum, doar am introdus variabilele în partea dreaptă și am apăsat butonul **OK**. Când utilizăm meniul, în fereastra care se deschide, observăm mai multe butoane. Cel care ne interesează aici este butonul **Statistics** (figura 6.1).

Acest submeniu ne permite să calculăm media, mediana, abaterea standard și diferite tipuri de percentile. Pe lângă acestea, putem alege să calculăm și alți indicatori ai tendinței centrale și variației cum ar fi modul, respectiv amplitudinea. De asemenea, în secțiunea **Distribution** putem calcula doi indicatori ai formei distribuției, **skewness** (alungirea) și **kurtosis** (aplatizarea), dar despre aceștia discutăm la secțiunea de explorare a datelor.

**Figura 6.1.** Meniul Frequencies, butonul Statistics

(a)



(b)

**Frequencies: Statistics**

**Percentile Values**

☐ Quantiles

☐ Cut p<sub>oints</sub> for: 10 equal groups

☒ P<sub>ercentile</sub>(s): 30

Add 10.0

Change 20.0

Remove

**Central Tendency**

☐ M<sub>ean</sub>

☐ M<sub>edian</sub>

☐ M<sub>ode</sub>

☐ S<sub>um</sub>

☐ Values are group midpoints

**Dispersion**

☐ S<sub>t</sub>d. d<sub>eviation</sub>

☐ V<sub>ariance</sub>

☐ R<sub>ange</sub>

☐ M<sub>inimum</sub>

☐ M<sub>aximum</sub>

☐ S. E. m<sub>ean</sub>

**Distribution**

☐ S<sub>kewness</sub>

☐ K<sub>urtosis</sub>

Continue Cancel Help

Fereastra afișată prin apăsarea butonului **Statistics** este intuitivă. Observăm că indicatorii sunt grupați în secțiunile **Percentile Values** (poziționare), **Central Tendency** (tendință centrală), **Dispersion** (variație) și **Distribution** (forma distribuției). În analiza noastră, suntem interesați să cunoaștem tendința centrală pentru fericire și sănătatea autoevaluată în rândul românilor. În baza de date WVS 2012, variabilele sunt V10 și V11. Indicatorii statistici sunt prezentați în tabelul 6.3a, iar tabelele de frecvențe sunt prezentate în tabelul 6.3b. Variabilele sunt ordinale: fericirea variază de la „deloc fericit” la „foarte fericit”, iar sănătatea autoevaluată variază de la „proastă” la „foarte bună”. Puteți schimba ordinea în care sunt așezate categoriile în funcție de codurile lor dacă, în meniul **Frequencies**, apăsați butonul **Format** și, în secțiunea **Order by**, bifați **Descending values**. Fiind variabile ordinale, putem calcula mediana și măsurile poziționării. În practică, deseori, întâlnim în multe lucrări și medii calculate pentru acest tip de variabilă ordinală.

**Tabelul 6.3.** Tabele de frecvență și indicatori statistici ai tendinței centrale și ai variației  
(a)

| Statistics     |         |                          |                                  |
|----------------|---------|--------------------------|----------------------------------|
|                |         | V10 Feeling of happiness | V11 State of health (subjective) |
| N              | Valid   | 1495                     | 1502                             |
|                | Missing | 8                        | 1                                |
| Mean           |         | 2.21                     | 2.28                             |
| Median         |         | 2.00                     | 2.00                             |
| Std. Deviation |         | .721                     | .830                             |
| Percentiles    | 25      | 2.00                     | 2.00                             |
|                | 50      | 2.00                     | 2.00                             |
|                | 75      | 3.00                     | 3.00                             |

(b)

| V10 Feeling of happiness |                    |           |         |               |                    |
|--------------------------|--------------------|-----------|---------|---------------|--------------------|
|                          |                    | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                    | 1 Very happy       | 205       | 13.6    | 13.7          | 13.7               |
|                          | 2 Rather happy     | 833       | 55.4    | 55.7          | 69.4               |
|                          | 3 Not very happy   | 397       | 26.4    | 26.6          | 96.0               |
|                          | 4 Not at all happy | 60        | 4.0     | 4.0           | 100.0              |
|                          | Total              | 1495      | 99.5    | 100.0         |                    |
| Missing                  | -2 No answer       | 4         | .3      |               |                    |
|                          | -1 Don't know      | 4         | .3      |               |                    |
|                          | Total              | 8         | .5      |               |                    |
| Total                    |                    | 1503      | 100.0   |               |                    |

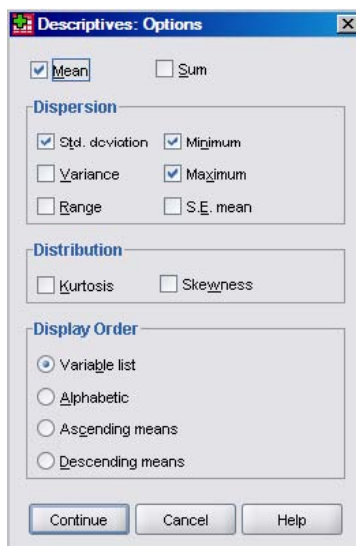
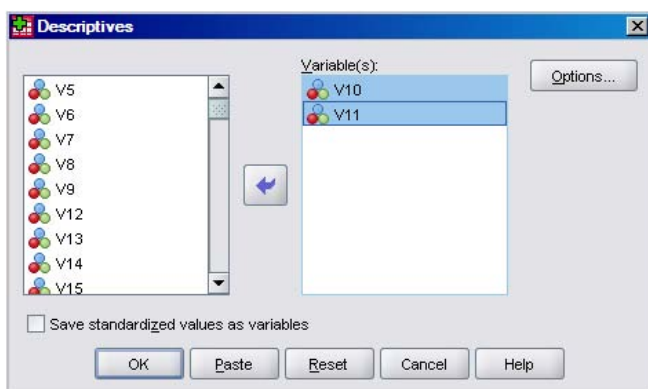
| V11 State of health (subjective) |              |           |         |               |                    |
|----------------------------------|--------------|-----------|---------|---------------|--------------------|
|                                  |              | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                            | 1 Very good  | 232       | 15.4    | 15.4          | 15.4               |
|                                  | 2 Good       | 747       | 49.7    | 49.7          | 65.2               |
|                                  | 3 Fair       | 390       | 25.9    | 25.9          | 91.1               |
|                                  | 4 Poor       | 134       | 8.9     | 8.9           | 100.0              |
|                                  | Total        | 1502      | 99.9    | 100.0         |                    |
| Missing                          | -2 No answer | 1         | .1      |               |                    |
| Total                            |              | 1503      | 100.0   |               |                    |

Mediana fericirii este egală cu 2, „destul de fericit”. Procentele ne arată că cel mai frecvent nivel de fericire ales de către români este „destul de fericit” (56%). Mediana stării de sănătate autoevaluate este egală cu 2, „bună”. Procentele ne arată că cea mai frecvent aleasă stare a sănătății de către români este „bună” (50%). Media este apropiată ca valoare de mediană pentru ambele variabile. Cuartilele

ne arată că cel puțin 25% dintre români au declarat că sunt „nu prea fericiți” sau „deloc fericiți” (percentila 75 = cuartila 3 = codul 3 „nu prea fericit”), respectiv că au o stare de sănătate „nu prea bună” sau chiar „proastă” (percentila 75 = cuartila 3 = codul 3 „nu prea bună”). Dacă dorim o informație mai detaliată, putem înlocui cuartilele cu decile, de exemplu : în secțiunea **Percentile Values** introducem valorile 10, 20, ... , 100 (figura 6.1b). Agresti și Finlay (2008) ne îndeamnă să fim precauți cu interpretarea atunci când variabila are puțin categorii (variante de răspuns).

Un alt meniu din care putem obține acești indicatori statistici este **Analyze > Descriptive Statistics > Descriptives**. Acesta este însă ceva mai limitat, permițând doar calcularea mediei și abaterii standard, fără mediană și percentile. Odată intrați în meniu (figura 6.2), apăsați butonul **Options** și alegem ce indicatori ne interesează. Rezultatele vor fi aceleași.

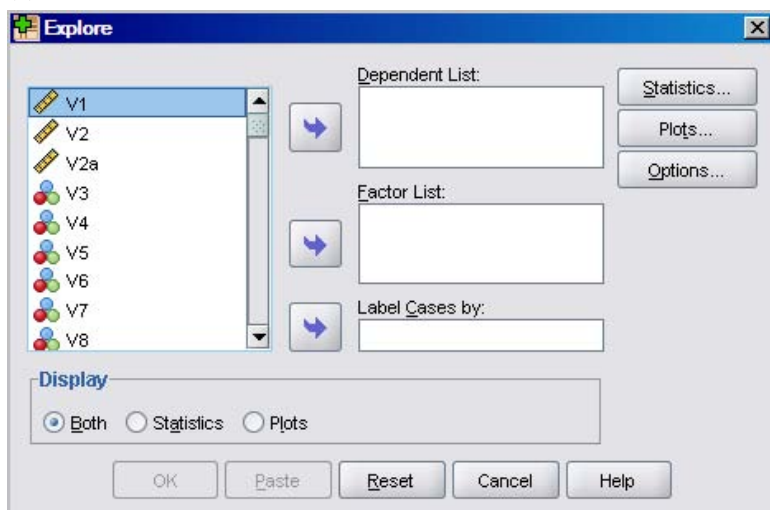
Figura 6.2. Meniul Descriptives



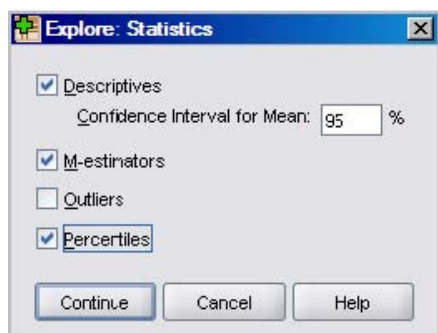
În fine, ultimul meniu prezentat aici, care poate fi folosit pentru calcularea acestor indicatori, este **Analyze > Descriptive Statistics > Explore** (figura 6.3a).

Figura 6.3. Meniul Explore

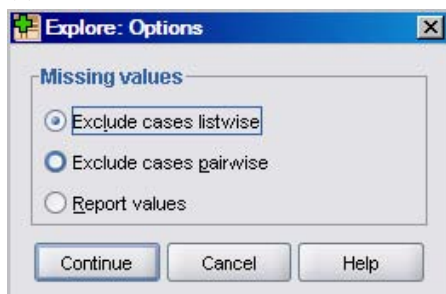
(a)



(b)



(c)





Acesta este ceva mai complex pentru că, așa cum ne arată și numele, este dedicat explorării datelor în vederea testării unor asumții de bază ale analizelor statistice uzual angajate în studiile sociale. Aici vom discuta doar despre cum obținem indicatorii discutați, restul meniului fiind abordat în secțiunea dedicată explorării datelor.

Fereastra are două secțiuni care ne interesează în acest moment : **Dependent List** și **Factor List**. La **Dependent List** introducem variabila pentru care dorim să calculăm statisticile, de exemplu, fericirea sau starea de sănătate autoevaluată. La **Factor List** introducem variabila care conține grupurile care urmează să fie comparate, de exemplu, mediul de rezidență, care conține două grupuri : locuitorii din urban și locuitorii din rural. Apăsând butonul **Statistics** putem selecta, pe lângă statisticile descriptive discutate, **M-estimators**, **Outliers** și **Percentiles** (figura 6.3b).

**M-estimators** sunt alternative robuste la medie și mediană. Bifând **Percentiles** obținem percentilele 5, 10, 25, 50, 75, 90 și 95. Prefer să lucrez cu meniul **Frequencies**, pentru că îmi dă mai multă libertate în opțiuni. Bifând **Outliers**, ne oferă un tabel cu ceea ce SPSS consideră a fi caz extrem (tabelul 6.4). Acest tabel nu este foarte informativ, pentru că oferă doar o selecție a așa-ziselor valori extreme. Coloana **Case Number** conține numărul rândului din **Data View**. Dacă am fi introdus în secțiunea **Label cases by** din fereastra principală (vezi figura 6.3a) o variabilă care conținea id-ul unic al fiecărui respondent, atunci tabelul ar mai fi conținut o coloană cu numele variabilei respective. Această alternativă este mai bună pentru că, dacă decidem să sortăm baza de date altfel decât în momentul în care am realizat tabelul (meniul **Sort Cases**), atunci informația din tabel devine inutilă.

**Tabelul 6.4.** Tabel Outliers obținut din meniul Explore

| Extreme Values   |         |   |             |                |
|--|---------|---|-------------|----------------|
|  |         |   | Case Number | Value          |
| V10 Feeling of happiness   | Highest | 1 | 621         | 4              |
|  |         | 2 | 622         | 4              |
|  |         | 3 | 623         | 4              |
|  |         | 4 | 624         | 4              |
|  |         | 5 | 625         | 4 <sup>a</sup> |
|  | Lowest  | 1 | 769         | 1              |
|  |         | 2 | 768         | 1              |
|  |         | 3 | 767         | 1              |
|  |         | 4 | 766         | 1              |
|  |         | 5 | 765         | 1 <sup>b</sup> |
| a. Only a partial list of cases with the value 4 are shown in the table of upper extremes. |         |   |             |                |
| b. Only a partial list of cases with the value 1 are shown in the table of lower extremes. |         |   |             |                |

În cadrul meniului **Explore**, o altă comandă care ne interesează acum este cea declanșată de butonul **Options** (figura 6.3c). Aici decidem cum sunt tratate nonrăspunsurile atunci când introducem, simultan, cel puțin două variabile la

**Dependent List** sau cel puțin două variabile la **Factor List**. Implicit, SPSS va trata nonrăspunsurile **listwise**, adică va dezactiva în analiză cazurile care au nonrăspunsuri. Decizia aparține însă cercetătorului.

Tabelul cu statistici oferit de meniul **Explore** conține multe informații utile (tabelul 6.5). Am calculat media, intervalul de încredere în jurul mediei, media calculată excluzând extremele distribuției (5% **Trimmed Mean**), mediana, varianța (pătratul abaterii standard), abaterea standard, valoarea minimă pe care o ia variabila, dar și valoarea maximă, amplitudinea (**range**), abaterea intercuartilă, alungirea (**skewness**) și aplătizarea (**kurtosis**).

**Tabelul 6.5.** Output produs de meniul Explore

| Descriptives                     |                                  |             |           |            |
|----------------------------------|----------------------------------|-------------|-----------|------------|
|                                  |                                  |             | Statistic | Std. Error |
| V10 Feeling of happiness         | Mean                             |             | 2.21      | .019       |
|                                  | 95% Confidence Interval for Mean | Lower Bound | 2.17      |            |
|                                  |                                  | Upper Bound | 2.25      |            |
|                                  | 5% Trimmed Mean                  |             | 2.19      |            |
|                                  | Median                           |             | 2.00      |            |
|                                  | Variance                         |             | .519      |            |
|                                  | Std. Deviation                   |             | .721      |            |
|                                  | Minimum                          |             | 1         |            |
|                                  | Maximum                          |             | 4         |            |
|                                  | Range                            |             | 3         |            |
|                                  | Interquartile Range              |             | 1         |            |
|                                  | Skewness                         |             | .305      | .063       |
|                                  | Kurtosis                         |             | .006      | .127       |
| V11 State of health (subjective) | Mean                             |             | 2.28      | .021       |
|                                  | 95% Confidence Interval for Mean | Lower Bound | 2.24      |            |
|                                  |                                  | Upper Bound | 2.32      |            |
|                                  | 5% Trimmed Mean                  |             | 2.26      |            |
|                                  | Median                           |             | 2.00      |            |
|                                  | Variance                         |             | .688      |            |
|                                  | Std. Deviation                   |             | .830      |            |
|                                  | Minimum                          |             | 1         |            |
|                                  | Maximum                          |             | 4         |            |
|                                  | Range                            |             | 3         |            |
|                                  | Interquartile Range              |             | 1         |            |
|                                  | Skewness                         |             | .365      | .063       |
|                                  | Kurtosis                         |             | -.336     | .127       |

Închei prin a atrage încă o dată atenția asupra stabilirii corecte a nivelului de măsurare al variabilei pentru care calculăm indicatorii statistici. Acest lucru se face înainte de realizarea calculelor respective. Deși aici am calculat media și

abaterea standard pentru variabile ordinale de tip Likert cu patru categorii, acest lucru nu înseamnă că acceptăm cu ușurință rezultatul primit. Vom întâlni în multe lucrări publicate astfel de analize. Trebuie să fim critici și să ne gândim cât de bine respectă cerințele de calcul astfel de măsurători și cât de interpretabil este rezultatul analizei.

## 6.2. Asocierea dintre variabile categoriale.

### Tabelul de contingență (Crosstabs)

După inspectarea individuală a variabilelor, vrem să vedem cum sunt asociate diferite variabile. De regulă, avem o variabilă a cărei variație dorim să o explicăm și mai multe variabile despre care credem că o influențează. Aici gândim bivariat. Cei din cuartila 1 de venit sunt mai mulțumiți cu viața lor decât cei din cuartila 2? Cei care au absolvit liceul sunt mai mulțumiți cu viața lor decât cei care au absolvit facultatea? Intuim deja de ce este util să învățăm și tehnici de analiză multivariată. Venitul mai mare crește posibilitatea de a satisface mai multe nevoi și aspirații, cum ar fi nevoia pentru o locuință cu mai multe camere, pentru o mașină mai încăpătoare, pentru vacanțe mai lungi etc. Cei care au absolvit niveluri formale de învățământ mai înalte au mai multe cunoștințe, lucru care le permite să fie mai flexibili pe piața muncii, să gestioneze riscurile mai ușor, să fie mai permeabili la schimbare etc. Însă, până la construirea unui model multivariat, nu putem face o idee despre obiectul studiului nostru folosind analizele bivariate. Decizia de a cumpăra un brand de cafea depinde de loialitatea față de brand? Dacă investigăm doar consumatori de cafea care nu sunt loiali nici unui brand, atunci când sunt la raftul de cafea, este culoarea ambalajului un factor de decizie pentru cumpărare?

Relația dintre două variabile categoriale poate fi observată folosind tabelul de contingență (**Crosstabs**). Variabilele categoriale sunt nominale sau ordinale. Esențial este ca, atunci când realizăm un tabel de contingență, ambele variabile să aibă puține categorii, pentru ca în fiecare celulă a tabelului să avem un număr rezonabil de cazuri. Un tabel cu 20 de rânduri și 10 coloane nu este util, pentru că, probabil, multe celule nu vor avea cazuri. Nu există o regulă care să specifice care este numărul optim de rânduri și coloane.

Persoanele care au încredere în semenii lor sunt mai fericite decât persoanele care nu au încredere în aceștia? Alți cercetători pot să pună întrebarea în sens invers: persoanele care sunt mai fericite au mai multă încredere în semenii lor decât persoanele care sunt mai puțin fericite? Sensul relației este stabilit printr-o atentă documentare teoretică. Programul de statistică nu alege variabila dependentă. El doar oferă calculele și graficele pe care le solicităm. Alegerea sensului relației este un act teoretic realizat înainte de a trece efectiv la analizarea datelor

în program. Domeniul fericirii este un exemplu foarte bun în ceea ce privește ambiguitatea direcției : de la fericire la altceva sau de la altceva la fericire. Pentru sociologi este specifică mai degrabă a doua variantă : presupunem că fericirea este starea la care trebuie să ajungem, trebuind să identificăm factorii care ne ajută în acest sens.

Să vedem care este relația dintre încredere și fericire. În WVS 2012, fericirea este măsurată prin întrebarea : „V10. Luând în considerare toate aspectele vieții dvs., ați spune că sunteți : 1. Foarte fericit ; 2. Destul de fericit ; 3. Nu prea fericit ; 4. Deloc fericit ?”. Încrederea este măsurată prin întrebarea : „V24. În general vorbind, ați spune că se poate avea încredere în cei mai mulți oameni sau că e mai bine să fii atent în relațiile cu oamenii : 1. Se poate avea încredere în cei mai mulți oameni ; 2. E mai bine să fii atent în relațiile cu oamenii ?”. Ne așteptăm că persoanele care au încredere în majoritatea oamenilor, adică aleg varianta 1 la V24, să fie mai fericite, adică aleg variantele 1 sau 2 la V10.

Avem două variabile categoriale : una nominală, V24, pe care o considerăm independentă, și una ordinală, V10, pe care o considerăm dependentă. Adică V10 este influențată de V24. Putem încrucișa aceste două variabile, pentru a vedea dacă presupunerea este corectă. Mai întâi realizăm câte un tabel de frecvență pentru V10 și V24, pentru (1) a vedea dacă există coduri de nonrăspuns care nu sunt declarate **missing** în program și pentru (2) a inspecta distribuția variabilelor. Dacă există coduri de nonrăspuns nedeclarate **missing**, atunci trebuie să mergem în **Variable View** > coloana **Missing** și să le declarăm. În ceea ce privește distribuția, ne interesează să avem suficiente cazuri pentru fiecare variantă de răspuns de la cele două variabile. S-ar putea ca la fericire, V10, să fie necesară o recodificare care presupune gruparea categoriilor. Există oameni care nu experimentează nici un pic de fericire (aleg varianta 4 la V10) ? În această situație s-ar putea să dorim unirea categoriilor „deloc fericit” și „nu prea fericit”. De asemenea, s-ar putea ca la încredere, V24, să nu avem variație, adică majoritatea să aibă sau să nu aibă încredere în semenii lor. În această situație, avem mai multe posibilități : considerăm că încrederea nu este măsurată bine și căutăm alt indicator pe care să îl folosim în analiză, sau considerăm că încrederea nu este un factor care afectează fericirea. Tabelele de frecvență sunt prezentate în tabelul 6.6.

Tabelul 6.6 prezintă informații despre cele două variabile. Respondenții au fost rugați să spună dacă cred că pot avea încredere în cei mai mulți dintre oameni sau e mai bine să ai grijă în relațiile cu oamenii. Majoritatea aleg a doua variantă de răspuns.

**Tabelul 6.6.** Tabele de frecvență : inspectarea variabilelor înainte de analiza de contingență (Crosstabs)

| <b>V10 Feeling of happiness</b> |                    |           |         |               |                    |
|---------------------------------|--------------------|-----------|---------|---------------|--------------------|
|                                 |                    | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                           | 1 Very happy       | 205       | 13.6    | 13.7          | 13.7               |
|                                 | 2 Rather happy     | 833       | 55.4    | 55.7          | 69.4               |
|                                 | 3 Not very happy   | 397       | 26.4    | 26.6          | 96.0               |
|                                 | 4 Not at all happy | 60        | 4.0     | 4.0           | 100.0              |
|                                 | Total              | 1495      | 99.5    | 100.0         |                    |
| Missing                         | -2 No answer       | 4         | .3      |               |                    |
|                                 | -1 Don't know      | 4         | .3      |               |                    |
|                                 | Total              | 8         | .5      |               |                    |
| Total                           |                    | 1503      | 100.0   |               |                    |

| <b>V24 Most people can be trusted</b> |                              |           |         |               |                    |
|---------------------------------------|------------------------------|-----------|---------|---------------|--------------------|
|                                       |                              | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                                 | 1 Most people can be trusted | 115       | 7.7     | 7.7           | 7.7                |
|                                       | 2 Need to be very careful    | 1373      | 91.4    | 92.3          | 100.0              |
|                                       | Total                        | 1489      | 99.0    | 100.0         |                    |
| Missing                               | -2 No answer                 | 4         | .2      |               |                    |
|                                       | -1 Don't know                | 11        | .7      |               |                    |
|                                       | Total                        | 14        | 1.0     |               |                    |
| Total                                 |                              | 1503      | 100.0   |               |                    |

Nonrăspunsurile sunt definite : codurile valide sunt grupate în rândul **Valid**, iar codurile de nonrăspuns sunt grupate în rândul **Missing**. Remarcăm, așa cum ne așteptam, că a patra categorie de fericire, „deloc fericit”, are o frecvență mult mai scăzută decât celelalte. Pentru moment, obiectivul nostru de cercetare este să vedem dacă încrederea este asociată cu fericirea sau nu. Este suficient, așadar, să am doar două categorii la variabila dependentă : fericiți și nefericiți. Așadar, folosind meniul **Transform > Recode into Different Variables**, vom crea o nouă variabilă dummy, cu numele v10rec, pornind de la V10 : codurile 1 și 2 devin 1, fericiți, iar codurile 3 și 4 devin 0, nefericiți. Mergând la încredere, observăm că majoritatea românilor nu au încredere în semenii lor, alegând varianta 2 de răspuns. Distribuția răspunsurilor ar putea proveni din modul în care este formulat itemul : nu reușește să discrimineze între indivizi. O discuție detaliată despre acest gen de situații poate fi consultată în Mărginean (1982). Pe de altă parte, aceasta ar putea fi realitatea în România anului 2012. Dacă ne uităm la distribuția răspunsurilor la această variabilă în alte țări incluse în studiu, vom observa că arată diferit : În Australia, 48% aleg a doua variantă, în Japonia, 61 %, în Noua Zeelandă, 43 %, în Suedia, 38% etc. Distribuții similare cu cea din țara noastră sunt întâlnite în Cipru, Peru etc. Acceptăm că putem folosi itemul în analiza noastră.

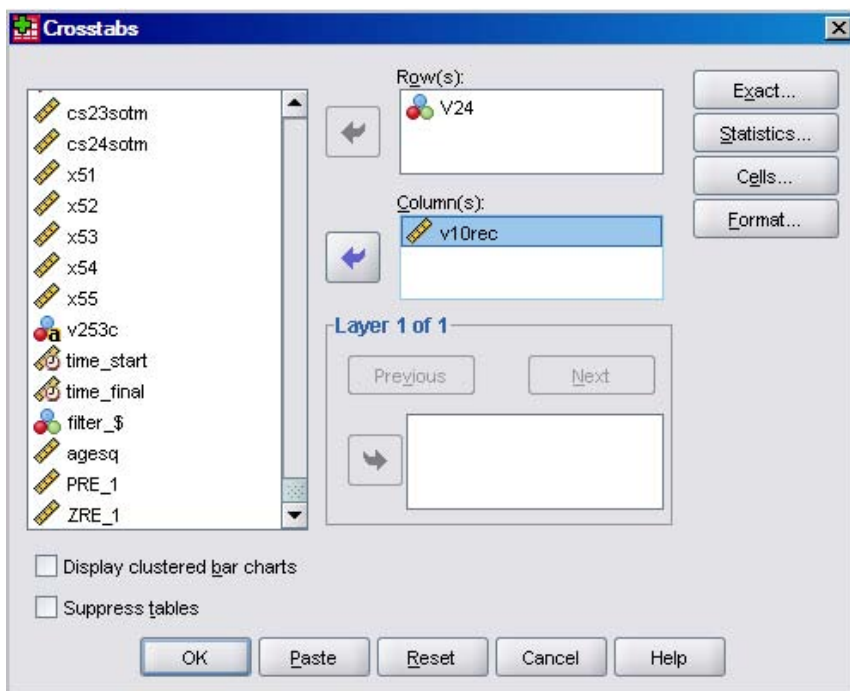
Tabelul de contingență este obținut din meniul **Analyze > Descriptive Statistics > Crosstabs**. Figura 6.4a prezintă fereastra principală care se deschide prin accesarea acestui meniu.

În stânga, observăm lista de variabile din care le alegem pe cele care ne interesează și le trecem în căsuțele din dreapta. La **Row(s)** introducem variabila care vrem să fie poziționată pe rândurile tabelului. La **Column(s)** introducem variabila care vrem să fie poziționată pe coloanele tabelului. Nu există o regulă cu privire la poziționarea pe rând sau coloană. Pe rând, e preferabil să introducem variabila cu cele mai multe categorii, iar pe coloană pe cea cu cele mai puține categorii. Astfel obținem un tabel care va fi mai ușor de încadrat într-o coală A4 orientată portret. În exemplul nostru, această discuție este irelevantă pentru că ambele variabile au doar două categorii de răspuns.

Dacă dorim să observăm relația dintre cele două variabile introduse în **Row(s)** și **Column(s)**, în funcție de valorile altei variabile, atunci vom utiliza **Layer 1 of 1**. De exemplu, vrem să vedem relația dintre încredere și fericire, în funcție de genul respondentului : care este relația pentru femei și care este relația pentru bărbați ? Folosind butonul **Next**, care se activează după ce introducem prima variabilă în **Layer 1 of 1**, putem subdivida și mai mult. Când folosim această opțiune, trebuie să avem destul de multe cazuri în eșantion pentru a fi relevante rezultatele.

Figura 6.4. Meniul Crosstabs

(a)



(b)

**Crosstabs: Cell Display**

**Counts**

☒ Observed  
☐ Expected

**Percentages**

☒ Row  
☐ Column  
☐ Total

**Residuals**

☐ Unstandardized  
☐ Standardized  
☒ Adjusted standardized

**Noninteger Weights**

☐ Round cell counts    ☐ Round case weights  
☐ Truncate cell counts    ☒ Truncate case weights  
☒ No adjustments

Continue Cancel Help

(c)

**Crosstabs: Statistics**

☒ Chi-square    ☐ Correlations

**Nominal**

☐ Contingency coefficient  
☐ Phi and Cramer's V  
☐ Lambda  
☐ Uncertainty coefficient

**Ordinal**

☐ Gamma  
☐ Somers' d  
☐ Kendall's tau-b  
☐ Kendall's tau-c

**Nominal by Interval**

☐ Eta

☐ Kappa  
☐ Risk  
☐ McNemar

☐ Cochran's and Mantel-Haenszel statistics  
Test common odds ratio equals: 1

Continue Cancel Help

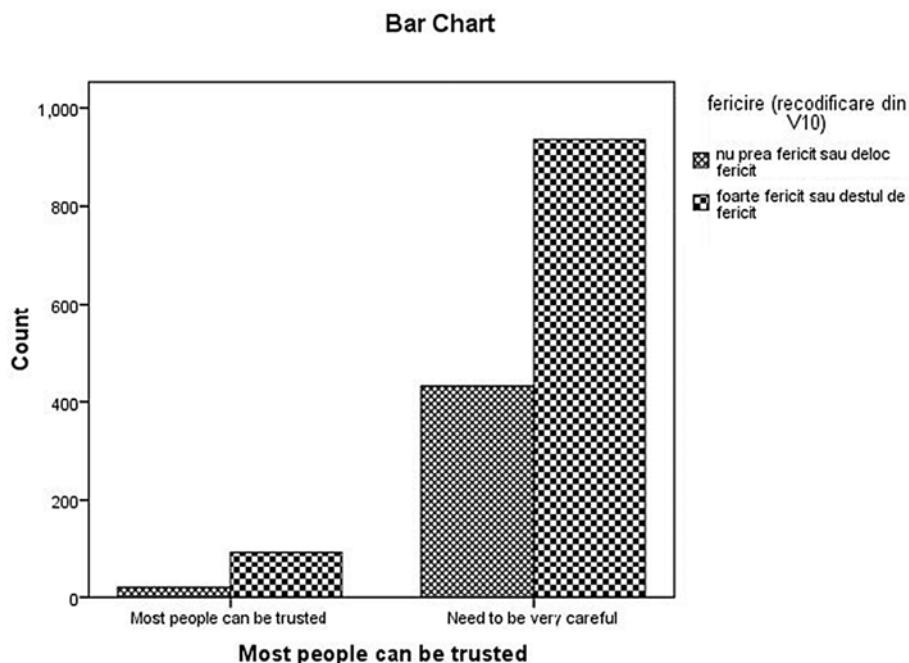
În colțul din stânga jos, observăm două opțiuni : **Display clustered bar charts** și **Suppress tables**. Prima produce un grafic bară similar cu cel din figura 6.5a. Forma prezentată aici este modificată față de cea produsă prin setările implicite de către SPSS. Modificările au fost făcute dând dublu-click pe graficul rezultat în **Output**. După ce s-a deschis pentru editare, selectăm pe rând fiecare bară. Apăsăm pe meniul **Edit > Properties** din fereastra **Chart Editor**. În fereastra **Properties** (figura 6.5b) aleg tabul **Fill & Border**. În secțiunea **Color** selectez căsuța colorată în alb și, la **Pattern**, modelul dorit. Repet operațiunea pentru cealaltă bară din tabel. Puteți realiza mai multe modificări din acest meniu, în funcție de preferințe.

Dacă bifăm cealaltă opțiune, **Suppress tables**, atunci rularea comenzii nu va afișa tabelul de contingență. Dacă ați bifat vreo opțiune în meniul care se deschide prin apăsarea butonului **Statistics**, atunci va fi afișat tabelul cu statisticile respective.

Scopul nostru principal este să vizualizăm sub formă de tabel relația dintre cele două variabile. Dacă după ce am introdus cele două variabile pe rând și pe coloană (figura 6.4a) apăsăm **OK**, tabelul rezultat va conține doar frecvențele absolute, adică numărul de persoane care au sau nu încredere în semenii lor și în starea acestora, de fericire sau nefericire (tabelul 6.7). 20 de persoane consideră că poți avea încredere în cei mai mulți oameni și se declară nefericiți. 94 de persoane consideră că se poate avea încredere în cei mai mulți oameni și se declară fericiți.

**Figura 6.5.** Grafic bară obținut folosind meniul Crosstabs

(a)





(b)



Tabelul 6.7. Tabel de contingență care conține doar frecvențe absolute (Count)

| V24 Most people can be trusted * v10rec fericire (recodificare din V10)<br>Crosstabulation |                                 |   |   |       |
|--|---------------------------------|---|---|-------|
| Count  |                                 |   |   |       |
|  |                                 | v10rec fericire (recodifi-<br>care din V10) |   | Total |
|  |                                 | 0 nu prea<br>fericit sau<br>deloc fericit   | 1 foarte<br>fericit sau<br>destul de<br>fericit |       |
| V24 Most people can<br>be trusted  | 1 Most people can be<br>trusted | 20  | 94  | 114   |
|  | 2 Need to be very<br>careful    | 432   | 935   | 1367  |
| Total  |                                 | 452   | 1029  | 1481  |

Este destul de greu să interpretăm datele vizualizate în acest mod. De aceea trebuie să transformăm frecvențele absolute în procente. Pentru că explicăm fericirea în funcție de încredere, vom calcula procente pe rând, pentru că pe rând am introdus variabila independentă. Apăsăm butonul **Cells**. Alegem, în secțiunea **Percentages**, căsuța **Row** (figura 6.4b). Astfel, totalul de 100% va fi pe fiecare rând. Tot aici am mai făcut o modificare față de setările implicite: în secțiunea **Noninteger Weights**, în loc de **Round cell counts**, am bifat **No adjustments**. Baza de date pe care sunt efectuate analizele aici este ponderată, iar ponderile au valori de tipul 1.410471 sau 0.780202. Dacă nu facem modificarea, atunci când calculează statisticile, programul va rotunji sau trunchia aceste valori. Rezultatul final nu va folosi ponderile în mod corespunzător.

Tabelul de contingență, care include frecvențele absolute (**Count**) și procente pe rând (**% within V24...**), este prezentat în tabelul 6.8. Pentru că am lăsat activă, în secțiunea **Counts**, opțiunea **Observed**, tabelul conține atât frecvențele absolute, cât și procente calculate din variabila încredere (totalurile pe rând sunt egale cu 100%). Pentru că am modificat opțiunea din secțiunea **Noninteger Weights**, frecvențele absolute au zecimale. În raport, folosim valorile rotunjite atât la frecvențele absolute, cât și la procente. În terminologia procentelor alese și calculate aici, 32% dintre cei care consideră că nu poți avea încredere în majoritatea oamenilor se declară nefericiți.

Tabelul este destul de dificil de citit, având prea multă informație. Este util să rămână vizibile doar procente, pentru a putea detecta dacă patternul așteptat prin ipoteza de lucru există sau nu. Putem să ne întoarcem în meniu și să debifăm opțiunea **Observed**, lăsând doar opțiunea **Row**. Dar mai rapid ar fi să edităm tabelul în **Output**.

**Tabelul 6.8.** Tabel de contingență care conține frecvențe absolute și procente pe rând

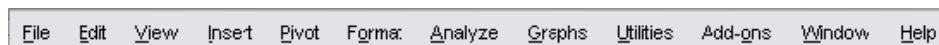
| <b>V24 Most people can be trusted * v10rec fericire (recodificare din V10)</b> |                                 |  |   |   |          |
|--|---------------------------------|--|---|---|----------|
| <b>Crosstabulation</b>   |                                 |  |   |   |          |
|  |                                 |  | v10rec fericire<br>(recodificare din V10) |   | Total    |
|  |                                 |  | 0 nu prea<br>fericit sau<br>deloc fericit | 1 foarte<br>fericit sau<br>destul de<br>fericit |          |
| V24 Most<br>people can<br>be trusted   | 1 Most people<br>can be trusted | Count                                      | 19.926                                    | 94.167  | 114.093  |
|  |                                 | % within V24 Most<br>people can be trusted | 17.5%                                     | 82.5%   | 100.0%   |
|  | 2 Need to be<br>very careful    | Count                                      | 431.848                                   | 935.166   | 1367.014 |
|  |                                 | % within V24 Most<br>people can be trusted | 31.6%                                     | 68.4%   | 100.0%   |
| Total  |                                 | Count                                      | 451.774                                   | 1029.333  | 1481.108 |
|  |                                 | % within V24 Most<br>people can be trusted | 30.5%                                     | 69.5%   | 100.0%   |

Așadar, în **Output**, dăm dublu click pe tabel. Tabelul se va deschide pentru editare. În acest moment, putem modifica etichetele, putem șterge sau adăuga

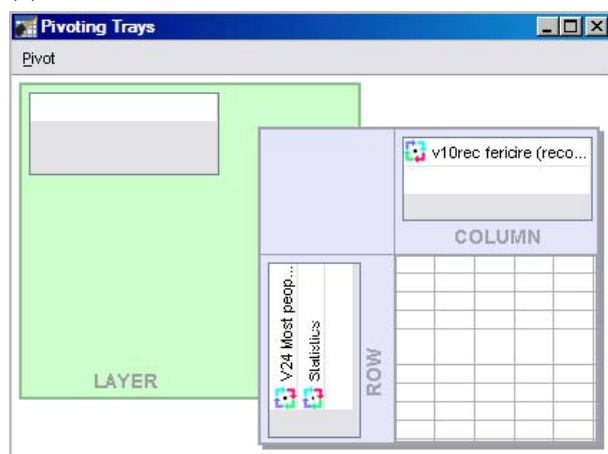
informație ș.a.m.d. Dar nu aceste lucruri ne interesează. Scopul nostru este să rămână vizibile doar procentele. Pentru aceasta, având tabelul deschis pentru editare, citim bara de meniuri și observăm că au apărut câteva opțiuni noi, printre care și **Pivot** (figura 6.6a). Meniul, înainte de dublu click, este identic cu cel din baza de date (**Data View** sau **Variable View**).

**Figura 6.6.** Editarea unui tabel de contingență în Output (Pivot)

(a)



(b)



(c)

### Crosstabs

[DataSet1] C:\Document

|                                |                            |  |  |          |
|--------------------------------|----------------------------|--|--|----------|
| V24 Most people can be trusted |                            | v10rec fericire (recodificare din V10) |  |          |
|                                |                            | 0 nu prea fericit sau deloc fericit    | 1 foarte fericit sau destul de fericit | Total    |
| V24 Most people can be trusted | Most people can be trusted | 19.926                                 | 94.167                                 | 114.093  |
|                                | 2 Need to be very careful  | 431.848                                | 935.166                                | 1367.014 |
| Total                          |                            | 451.774                                | 1029.333                               | 1481.108 |

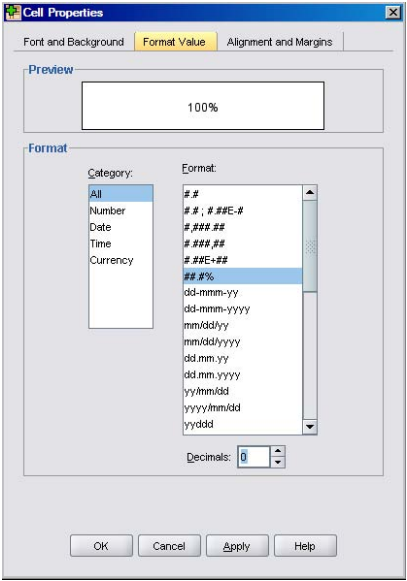
(d)

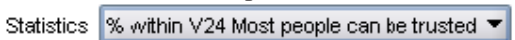
| V24 Most people can be trusted * v10rec fericire (recodificare din V10)<br>Crosstabulation |                              |  |  |        |
|--|------------------------------|--|--|--------|
| Statistics=% within V24 Most people can be trusted   |                              |  |  |        |
|  |                              | v10rec fericire (recodificare din V10) |  | Total  |
|  |                              | 0 nu prea fericit sau deloc fericit    | 1 foarte fericit sau destul de fericit |        |
| V24 Most people can be trusted   | 1 Most people can be trusted | 17.5%                                  | 82.5%                                  | 100.0% |
|  | 2 Need to be very careful    | 31.6%                                  | 68.4%                                  | 100.0% |
| Total  |                              | 30.5%                                  | 69.5%                                  | 100.0% |

(e)

| V24 Most people can be trusted * v10rec fericire (recodificare din V10) Crosstabulation |                              |  |  |        |
|---|------------------------------|--|--|--------|
| Statistics % within V24 Most people can be trusted ▼                                    |                              |  |  |        |
|   |                              | v10rec fericire (recodificare din V10) |  | Total  |
|   |                              | 0 nu prea fericit sau deloc fericit    | 1 foarte fericit sau destul de fericit |        |
| V24 Most people can be trusted  | 1 Most people can be trusted | 17.5%                                  | 82.5%                                  | 100.0% |
|   | 2 Need to be very careful    | 31.6%                                  | 68.4%                                  | 100.0% |
| Total   |                              | 30.5%                                  | 69.5%                                  | 100.0% |

(f)



La **Pivot**, selectăm **Pivoting Trays** și se deschide fereastra din figura 6.6b. Fereastra aceasta are două elemente: tabelul (fereastra din plan apropiat) și layerul (fereastra din plan îndepărtat). În tabel, pe coloană (**COLUMN**) avem variabila **v10rec**, iar pe rând (**ROW**) avem variabila **V24** și statisticile calculate, mai exact frecvențele absolute și procentele pe rând. Ducem mouse-ul pe textul **Statistics** din **ROW** și, ținând apăsat, folosind procedeul drag-and-drop, tragem de acesta până în colțul din stânga sus de la **Layer**, care este alb. Când ajungem pe suprafața albă, săgeata mouse-ului se va transforma într-o mână. În acest moment eliberăm **Statistics**. În **Output**, observăm cum s-a modificat tabelul (figura 6.6c). În acest moment afișează doar frecvențele absolute (**Count**). Pentru că suntem interesați să afișeze procentele, ducem mouse-ul pe butonul de deasupra tabelului  și alegem **% within V24**. Modificarea este instantanee, putând citi imediat informațiile (figura 6.6d). Pentru a închide fereastra de editare a tabelului, este suficient ca, în **Output**, să dăm click în afara lui.

Pentru că nu raportăm procentele cu virgulă, vom da iarăși dublu click pe tabel. Selectăm toate celulele cu procente în ele (figura 6.6e) și, în meniul care se deschide, vom selecta **Format > Cell Properties > Format value > Decimals = 0** (figura 6.6e). Apăsăm butonul **Apply** și apoi **OK**.

Ipoteza spune că oamenii care au încredere în semenii lor sunt mai fericiți: 83% dintre cei care au încredere în semenii lor sunt fericiți și 68% dintre cei care nu au încredere în semenii lor sunt fericiți. Procentele par să susțină ideea noastră.

Înainte însă ar fi util să rulăm un test de semnificație. Discuția în detaliu, despre ce sunt testele de semnificație, care sunt argumentele pro și contra utilizării lor ș.a., depășește scopul acestei lucrări. Cititorul este rugat să consulte lucrările dedicate acestui subiect, având în vedere importanța pe care o au în analizele statistice. De asemenea, este rugat să înțeleagă care este relația cu utilizarea intervalelor de încredere pentru realizarea de inferențe. Pentru înțelegerea corectă a acestui concept, trebuie să înțeleagă diferența între populație și eșantion, parametru și statistică, eșantion probabilist și eșantion neprobabilist, trebuie înțelese concepte precum probabilitate, distribuție de eșantionare etc. Utile în acest sens sunt lucrările scrise de Henkel (1976) și de Mohr (1990). Din ipoteza de cercetare, sunt derivate o serie de ipoteze statistice. Un test de semnificație caută să verifice dacă putem să respingem ipoteza de nul. Ipoteza de nul, așa cum sugerează numele acesteia, presupune, de exemplu, că două variabile sunt independente, adică nu au nici o relație, nu sunt asociate. Cercetătorul testează această ipoteză folosind un eșantion probabilist extras din populația de referință pentru studiul său. De exemplu, testăm independența dintre încredere și fericire folosind un eșantion reprezentativ pentru populația care are 18 ani sau peste, locuiește în România, nefiind instituționalizată. Pornind de la acest eșantion care are, să zicem,

un volum de 1.500 de persoane, cercetătorul va face inferențe pentru întreaga populație din care a fost extras. Acesta este însă unul dintre eșantioanele care puteau fi extrase folosind aceeași schemă de eșantionare. Dacă aplicăm aceiași pași și același algoritm, vor rezulta eșantioane care includ alte persoane decât eșantioanele extrase anterior. În ce măsură rezultatul din eșantionul nostru se datorează întâmplării? Există o relație între încredere și fericire în cadrul populației?

Pentru a răspunde la această întrebare, putem folosi testul Pearson chi-square sau, dacă dorim să îl citim în limba română, hi pătrat. Acest test presupune realizarea unui tabel de contingență: sunt comparate frecvențele observate din fiecare celulă cu frecvențele așteptate din pură întâmplare pentru celulele respective. Calculele sunt explicate detaliat în multe lucrări de statistică, cum ar fi cea scrisă de Field (2009). Înainte de a calcula acest test, trebuie să alegem o valoare teoretică standard a nivelului de semnificație cu care să o comparăm pe cea calculată de program. În științele sociale, cele mai utilizate sunt 0.05, pentru un nivel de încredere de 95%, și 0.01, pentru un nivel de încredere de 99%. Dacă alegem valoarea teoretică 0.05, iar cea calculată este mai mică decât aceasta, să zicem 0.02, atunci putem respinge ipoteza de nul a independenței celor două variabile. Probabilitatea de a greși spunând că încrederea și fericirea sunt asociate este mică. Atenție la limbaj: discuția se poartă în termeni probabilistici. Nu putem spune: „sigur există o relație”, ci „este mai probabil să existe decât să nu existe”. În practică, când  $p$  calculat de SPSS este mai mic decât pragul teoretic utilizat, să zicem 0.05, spunem că relația este semnificativă statistic. Sau, și mai scurt, că relația este semnificativă. Să nu confundăm însă sensul de aici cu ideea de relație puternică. Înseamnă doar că putem respinge ipoteza de nul, nu și că relația este puternică. Înseamnă doar că probabilitatea de a greși spunând că fericirea este asociată cu încrederea este mai mică de 0.05 sau 5%. Aceste praguri teoretice sunt relativ arbitrare. Nu există o justificare solidă teoretic pentru alegerea lor. De ce un  $p$  calculat egal cu 0.06 face relația nesemnificativă statistic, iar un  $p$  calculat egal cu 0.05 o face semnificativă statistic? S-a dezvoltat o literatură alternativă pe acest subiect care merită consultată (Kline, 2004).

Testul chi-square este obținut apăsând butonul **Statistics**: în fereastra care se deschide, bifăm **Chi-square** (figura 6.4c). Rezultatele pentru analiza asocierii dintre încredere și fericire sunt prezentate în tabelul 6.9. Ne interesează primul rând. Coloana **Asymp. Sig. (2-sided)** conține valoarea  $p$  calculată. Aici este egală cu 0.002. O comparăm cu valoarea teoretică 0.05, aleasă înainte de a rula analiza. În sine, testul chi-square nu ne spune mare lucru și, dacă nu ținem cont de anumite asumptii ale acestuia, poate chiar să dezinformeze (Reynolds, 1984). Aici ne spune că încrederea și fericirea sunt asociate statistic:  $p$  calculat = 0.002, valoare mai mică decât 0.05. Acest rezultat are un grad de acuratețe ridicat dacă celulele tabelului de contingență conțin o anumită frecvență așteptată (vezi prima notă de sub tabel: **0 cells (0.0%) have expected count less than 5**). În eșantioanele cu multe unități, este foarte posibil ca  $p$  calculat să fie mai mic decât 0.05, chiar dacă variabilele sunt slab asociate. Argumentele teoretice pentru investigarea

acestei relații trebuie să fie bine gândite. O metodă empirică de verificare a acestei situații presupune utilizarea unei alte informații pe care ne-o poate calcula SPSS : calcularea valorilor reziduale ajustate (**adjusted standardized residuals**) (figura 6.4b). Acestea ne arată care celule explică asocierea dintre cele două variabile (p calculat la chi-square mai mic decât 0.05). Într-un tabel 2x2, adică dintre două variabile dihotomice, nu este prea relevant să ne uităm la aceste statistici, însă într-unul care are cel puțin o variabilă cu mai mult de două categorii se pot dovedi foarte utile în explicație (Field, 2009).

**Tabelul 6.9.** Testul Pearson chi-square : valoare și p

| Chi-Square Tests  |                    |    |                       |                      |                      |
|---|--------------------|----|-----------------------|----------------------|----------------------|
|   | Value              | df | Asymp. Sig. (2-sided) | Exact Sig. (2-sided) | Exact Sig. (1-sided) |
| Pearson Chi-Square  | 9.912 <sup>a</sup> | 1  | .002                  |                      |                      |
| Continuity Correction <sup>b</sup>  | 9.257              | 1  | .002                  |                      |                      |
| Likelihood Ratio  | 10.917             | 1  | .001                  |                      |                      |
| Fisher's Exact Test   |                    |    |                       | .001                 | .001                 |
| Linear-by-Linear Association  | 9.905              | 1  | .002                  |                      |                      |
| N of Valid Cases  | 1481               |    |                       |                      |                      |
| a. 0 cells (0.0%) have expected count less than 5. The minimum expected count is 34.80. |                    |    |                       |                      |                      |
| b. Computed only for a 2x2 table  |                    |    |                       |                      |                      |

În literatura de specialitate, starea civilă este considerată un predictor al fericirii. Unii autori susțin că implicarea într-o relație de cuplu crește fericirea partenerilor (Zimmerman și Easterlin, 2006). Folosind datele WVS 2012, putem inspecta, într-o primă fază, relația dintre fericire și starea civilă. Folosesc fericirea recodificată similar cu exemplul discutat anterior. Starea civilă are trei categorii : căsătorit sau angajat într-o relație de cuplu ; divorțat, separat sau singur ; văduv. Cele două variabile se numesc v10rec, respectiv v57rec. Rezultatul, incluzând reziduurile ajustate, este prezentat în tabelul 6.10. Valoarea lui chi-square este 90.692 (2 grade de libertate), iar valoarea lui p calculat este mai mică decât 0.01. Relația dintre starea civilă și fericire este probabilă. Mai mult, dacă ne uităm la reziduurile ajustate, aceasta este dată de fiecare tip de stare civilă. Reziduurile ajustate mai mari de 1.96, ignorând semnul, arată o relație semnificativă la nivel de celulă pentru un nivel de încredere de 95 %. Reziduurile ajustate mai mari de 2.58, ignorând semnul, arată o relație semnificativă la nivel de celulă pentru un nivel de încredere de 99 %. Semnele ne arată direcția relației. Reziduul -4.4 ne arată că persoanele care au o relație de cuplu, formalizată sau nu, nu trăiesc o stare de nefericire. În schimb, reziduul 9.5 ne arată că persoanele văduve trăiesc o stare de nefericire. Decesul partenerului reprezintă o pierdere grea în viața unei persoane. Trebuie investigat mai departe care sunt factorii care îi ajută pe cei divorțați, separați sau singuri să compenseze efectul pozitiv al prezenței unui partener de viață. Analizele multivariate se dovedesc utile în acest sens.

**Tabelul 6.10.** Tabel de contingență cu statistici : stare civilă și fericire

| <b>v57rec starea civila (recodificare din V57) * v10rec fericire (recodificare din V10) Crosstabulation</b> |                                      |  |  |   |       |
|---|--------------------------------------|--|--|---|-------|
|   |                                      | Statistics   | v10rec fericire<br>(recodificare din V10)    |   | Total |
|   |                                      |  | 0 nu prea<br>fericit sau<br>deloc<br>fericit | 1 foarte<br>fericit sau<br>destul de<br>fericit |       |
| v57rec starea<br>civila (recodifi-<br>care din V57)   | 1 casatorit<br>sau are o<br>relatie  | Count  | 263  | 720   | 983   |
|   |                                      | % within v57rec<br>starea civila (recodifi-<br>care din V57) | 27%  | 73%   | 100%  |
|   |                                      | Adjusted Residual  | -4.4   | 4.4   |       |
|   | 2 divortat,<br>separat sau<br>singur | Count  | 91   | 258   | 348   |
|   |                                      | % within v57rec<br>starea civila (recodifi-<br>care din V57) | 26%  | 74%   | 100%  |
|   |                                      | Adjusted Residual  | -2.1   | 2.1   |       |
|   | 3 vaduv                              | Count  | 100  | 58  | 158   |
|   |                                      | % within v57rec<br>starea civila (recodifi-<br>care din V57) | 63%  | 37%   | 100%  |
|   |                                      | Adjusted Residual  | 9.5  | -9.5  |       |
| Total   |                                      | Count  | 454  | 1036  | 1490  |
|   |                                      | % within v57rec<br>starea civila (recodifi-<br>care din V57) | 30%  | 70%   | 100%  |

Într-o altă analiză, am putea fi interesați să vedem dacă există o asocierie între starea civilă și încrederea în oameni. Așteptarea noastră este că persoanele care au ieșit dintr-o relație de cuplu vor fi mai reticente în a se încrede în alte persoane. Tabelul 6.11 prezintă rezultatele analizei bivariate. Valoarea lui chi-square este 9.680 (2 grade de libertate), iar valoarea p calculată este mai mică decât 0.01. Putem respinge ipoteza de nul a independenței celor două variabile. Inspectând reziduurile ajustate, observăm că relația dintre cele două variabile se datorează în principal statutului de divorțat, separat sau singur, deoarece reziduurile ajustate din dreptul acestei categorii au valori mai mari de 2.58, ignorând semnul, pe când celelalte au valori mai mici de 1.96, ignorând semnul. Așteptarea noastră este însă confirmată parțial, deoarece aceștia consideră că se poate avea încredere în majoritatea oamenilor (reziduul ajustat = 3.0). Trebuie investigate motivele pentru care se întâmplă acest lucru. Ipoteza noastră se baza pe ideea că relația de cuplu încetează pentru că cel puțin unul dintre parteneri a găsit o alternativă mai bună din diferite puncte de vedere. Astfel, celălalt se va simți trădat. Însă datele ne largesc orizontul. Controlând și pentru această idee, ar trebui văzut în ce măsură persoanele care



aleg să caute un alt parteneriat sunt mai deschise la experimentare, mai permeabile la schimbare etc. Analiza multivariată se dovedește din nou utilă. Trebuie să ne întoarcem la teorie, să o analizăm mai atent, să vedem ce ne-a scăpat și să construim un model explicativ pe care să îl testăm folosind o analiză care permite utilizarea simultană a mai multor variabile independente.

**Tabelul 6.11.** Tabel de contingență cu statistici : stare civilă și încredere în oameni

| <b>v57rec starea civila (recodificare din V57) * V24 Most people can be trusted</b> |                                |  |                                |                           |       |
|---|--------------------------------|--|--------------------------------|---------------------------|-------|
| <b>Crosstabulation</b>  |                                |  |                                |                           |       |
|   |                                | Statistics   | V24 Most people can be trusted |                           | Total |
|   |                                |  | 1 Most people can be trusted   | 2 Need to be very careful |       |
| v57rec starea civila (recodificare din V57)   | 1 casatorit sau are o relatie  | Count  | 67                             | 914                       | 980   |
|   |                                | % within v57rec starea civila (recodificare din V57) | 7%                             | 93%                       | 100%  |
|   |                                | Adjusted Residual                                    | -1.8                           | 1.8                       |       |
|   | 2 divortat, separat sau singur | Count  | 40                             | 309                       | 349   |
|   |                                | % within v57rec starea civila (recodificare din V57) | 11%                            | 89%                       | 100%  |
|   |                                | Adjusted Residual                                    | 3.0                            | -3.0                      |       |
|   | 3 vaduv                        | Count  | 8                              | 147                       | 154   |
|   |                                | % within v57rec starea civila (recodificare din V57) | 5%                             | 95%                       | 100%  |
|   |                                | Adjusted Residual                                    | -1.3                           | 1.3                       |       |
| Total   |                                | Count  | 114                            | 1369                      | 1483  |
|   |                                | % within v57rec starea civila (recodificare din V57) | 8%                             | 92%                       | 100%  |

În meniul din figura 6.4c, putem alege dintre mai mulți indicatori de asociere și chiar de corelație. Aceștia sunt grupați în funcție de tipul variabilelor pe care dorim să le asociem : nominale cu nominale (**Contingency coefficient, Phi and Cramer's V, Lambda, Uncertainty coefficient**), ordinale cu ordinale (**Gamma, Somers' d, Kendall's tau-b, Kendall's tau-c**), metrice cu metrice (**Correlations**) etc. Logica acestora va fi înțeleasă după parcurgerea capitolului dedicat corelației metrice. Spre deosebire de testul chi-square, aceștia sunt indicatori care iau o valoare într-un interval și ne arată direcția și intensitatea relației. O prezentare excelentă a diferenței dintre aceștia și a momentelor în care este potrivit să îl utilizăm pe unul sau altul a fost realizată de Chen și Popovich (2002).

### 6.3. Diferențe între medii : testul t pentru eșantioane independente și ANOVA

Uneori suntem interesați să vedem dacă două grupuri au valori similare pentru o anumită caracteristică. În situații puțin mai complexe, am putea fi interesați să comparăm trei sau mai multe grupuri după o anumită caracteristică. Lucrăm cu două variabile simultan : una categorială, care dă grupurile, și una metrică, pentru care calculăm media în cadrul fiecărui grup. Variabila categorială este cea care explică – variabilă independentă. Variabila metrică este cea explicată – variabilă dependentă. Într-un studiu care analizează discriminarea femeilor pe piața muncii, putem fi interesați să comparăm salariul lunar al femeilor și bărbaților care au locuri de muncă similare. Variabila de grupare va fi sexul, iar cea pentru care calculăm mediile va fi salariul lunar. Într-o cercetare de marketing, putem fi interesați să comparăm volumul vânzărilor, într-o anumită perioadă, pentru anvelopele de iarnă și pentru cele de vară produse de o anumită companie. Variabila de grupare va fi tipul de anvelope, iar cea pentru care calculăm mediile va fi volumul vânzărilor.

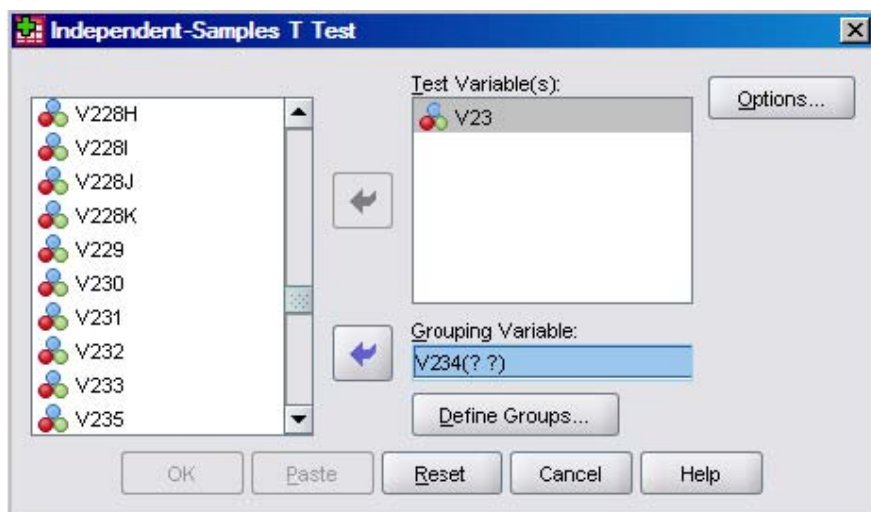
Ipoteza de nul va fi că mediile grupurilor comparate sunt egale. Dacă valoarea p calculată este mai mică decât pragul critic ales ca referință, 0.05 sau 0.01, atunci respingem ipoteza de nul și considerăm plauzibilă ipoteza alternativă. Grupurile comparate diferă în ceea ce privește caracteristica respectivă. Salariul bărbaților ar putea fi mai mare decât cel al femeilor. Trebuie investigate motivele acestei situații. Anvelopele de vară ale companiei sunt vândute într-o cantitate mai mare decât anvelopele de iarnă. Trebuie aflat de ce se întâmplă acest lucru. Observăm că testul t pentru eșantioane independente sau analiza de varianță (ANOVA) ne deschid căi interesante pentru explicarea unei situații. Însă, de regulă, cercetătorul nu se limitează la ele, ci, folosind modele explicative, aplică diferite tehnici de analiză multivariată pentru a reprezenta cât mai adecvat realitatea socială, ceva mai complexă decât aceste relații bivariate. ANOVA este necesară pentru că, dacă avem cel puțin trei grupuri și aplicăm câte un test t în cazul fiecărei perechi, există posibilitatea să vedem diferențe chiar și acolo unde nu există. Adică respingem ipoteza de nul când nu trebuie (Henkel, 1976). Testele de semnificații despre care discutăm ne arată dacă diferențele dintre mediile grupurilor există datorită variațiilor aleatoare de la un eșantion la altul ori pentru că datele provin din populații în care mediile chiar sunt diferite (Iversen și Norpoth, 1987). Pentru a fi relevantă comparația, grupurile trebuie să difere doar în ceea ce privește caracteristica presupusă a da diferența. Trebuie să aibă variații similare (Iversen și Norpoth, 1987). Pentru verificarea acestei asumptii, există mai multe teste de semnificații. SPSS ne oferă testul Levene. În funcție de informația arătată de acesta, interpretăm și rezultatul testului t pentru eșantioane independente și ANOVA. O altă condiție este ca variabila metrică să fie distribuită normal. Dacă

grupurile comparate au mărimi și variații diferite, iar distribuția este alungită sever, atunci este destul de probabil ca rezultatul analizelor să nu fie adecvat (Agresti și Finlay, 2008). Atunci când aceste asumptii nu pot fi satisfăcute, ar fi util să înlocuim sau măcar să comparăm rezultatele celor două analize cu cele ale echivalentelor nonparametrice care pot fi calculate în SPSS. Dar aceasta este o altă discuție.

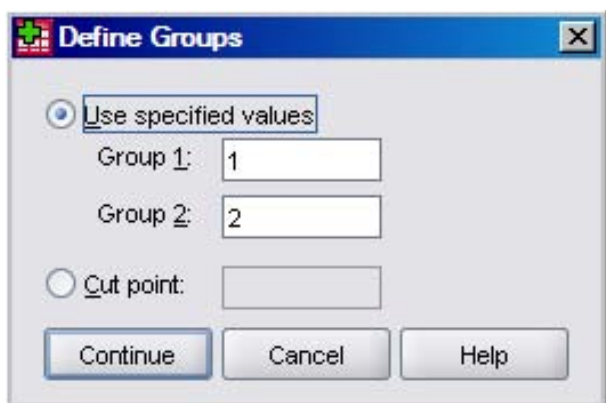
Testul t pentru eșantioane independente poate fi calculat folosind meniul **Analyze > Compare Means > Independent-Samples T Test** (figura 6.7). Fereastra care se deschide este foarte intuitivă pentru utilizator: în partea stângă avem lista de variabile din care le alegem pe cele pe care dorim să le utilizăm în analiză. În **Test Variable(s)** introducem variabila metrică pentru care dorim să calculăm media. Aici se introduce satisfacția cu viața care are numele V23 în WVS 2012. În **Grouping Variable**, introducem variabila categorială care dă grupurile comparate după variabila metrică. Aici se introduce variabila V234: „Slujba dvs. presupune să aveți/să fi avut pe cineva în subordine? 1. Da, 2. Nu”. În figura 6.7a observăm că în dreptul numelui variabilei sunt, între paranteze, două semne de întrebare: **V234(? ?)**. SPSS solicită codurile celor două grupuri pentru care dorim să comparăm mediile satisfacției cu viața. Le aflăm dintr-un tabel de frecvență. Aici corespund chestionarului: 1 înseamnă că respondentul are persoane în subordine la locul de muncă, iar 2 că nu are. Pentru a le introduce, apăsăm butonul **Define Groups** (figura 6.7b). Pentru că știm exact ce grupuri dorim să comparăm, variabila având oricum doar două coduri valide, selectăm **Use specified values** și, la **Group 1**, respectiv **Group 2**, introducem cele două coduri (figura 6.7b). Apăsăm **Continue**, iar rezultatul este vizibil în figura 6.7c.

Figura 6.7. Meniul Independent-Samples T Test

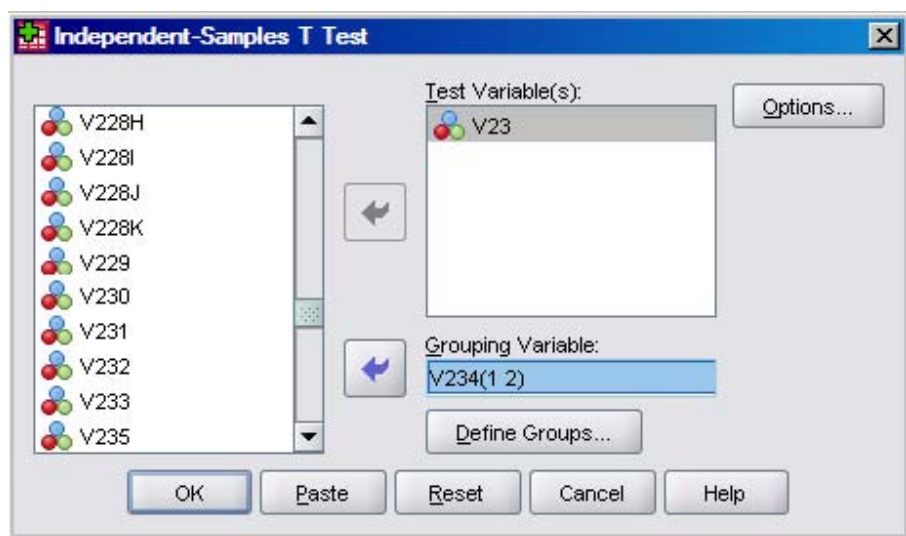
(a)



(b)



(c)



Rezultatul analizei este prezentat în tabelul 6.12. Mai întâi, sunt afișate câteva statistici descriptive. Sunt 358 de persoane care au subordonați, spre deosebire de cei care nu au alte persoane în subordine, care sunt în număr de 794. Cei din urmă au o medie a satisfacției cu viața egală cu 6.94, iar cei din urmă egală cu 6.63. Abaterile standard sunt apropiate ca valoare, 2.28 și 2.35. Satisfacția cu viața este măsurată pe o scală de la 1 la 10, scorurile mari indicând o satisfacție cu viața mai ridicată. Următoarea figură conține testul Levene și testul t pentru eșantioane independente. Testul Levene ne spune că varianțele celor două grupuri sunt egale. Valoarea p calculată pentru acesta este egală cu 0.112. Fiind mai mare decât pragul teoretic de 0.05, nu putem respinge ipoteza de nul a egalității

variantelor. Din acest motiv, o să citim testul t de pe rândul **Equal variances assumed**. Dacă valoarea p calculată a testului Levene ar fi fost mai mică decât 0.05, atunci am fi citit testul t de pe rândul **Equal variances not assumed**. Testul t ne spune că cele două medii sunt diferite : valoarea p calculată este egală cu 0.038, care este mai mică decât pragul teoretic de 0.05. Putem respinge ipoteza de nul a similarității satisfacției cu viața în rândul celor două grupuri : cu sau fără subordonați la locul de muncă. Din punct de vedere statistic, rezultatul ar putea fi satisfăcător. Totuși, cercetătorul nu trebuie să se mulțumească cu o abordare empiristă a realității sociale. Cele două medii sunt diferite prin 0.3 unități pe o scală de 10 puncte. Este aceasta o diferență de luat în seamă din punct de vedere practic ?

**Tabelul 6.12.** Testul t pentru eșantioane independente : output

| Group Statistics                |                                  |     |      |                |                 |
|---------------------------------|----------------------------------|-----|------|----------------|-----------------|
|                                 | V234 Are you supervising someone | N   | Mean | Std. Deviation | Std. Error Mean |
| V23 Satisfaction with your life | 1 yes                            | 358 | 6.94 | 2.281          | .121            |
|                                 | 2 no                             | 794 | 6.63 | 2.359          | .084            |

| Independent Samples Test        |                             |   |      |                              |         |                 |                 |                       |   |       |
|---------------------------------|-----------------------------|---|------|------------------------------|---------|-----------------|-----------------|-----------------------|---|-------|
|                                 |                             | Levene's Test for Equality of Variances |      | t-test for Equality of Means |         |                 |                 |                       |   |       |
|                                 |                             |   |      |                              |         |                 |                 |                       | 95% Confidence Interval of the Difference |       |
|                                 |                             | F                                       | Sig. | t                            | df      | Sig. (2-tailed) | Mean Difference | Std. Error Difference | Lower                                     | Upper |
| V23 Satisfaction with your life | Equal variances assumed     | 2.531                                   | .112 | 2.077                        | 1150    | .038            | .309            | .149                  | .017                                      | .600  |
|                                 | Equal variances not assumed |   |      | 2.104                        | 710.347 | .036            | .309            | .147                  | .021                                      | .597  |

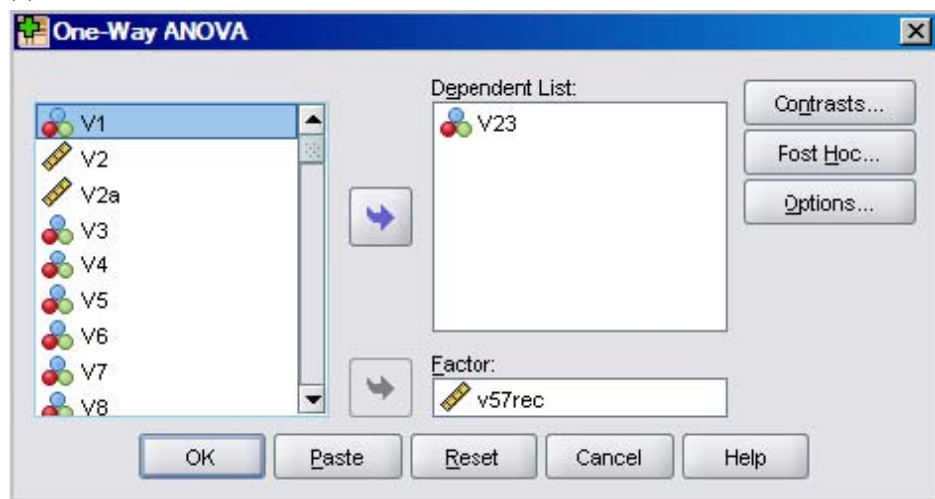
Atunci când avem mai multe grupuri ce trebuie comparate, utilizăm analiza de varianță, prescurtată ANOVA. Nivelul satisfacției cu viața diferă în funcție de starea civilă ? Am grupat persoanele în trei categorii : 1. Căsătorit sau trăiesc împreună cu cineva, dar nu suntem căsătoriți ; 2. Divorțat, separat (despărțit

nelegal) sau necăsătorit și fără a locui cu un partener ; 3. Văduv. Ipoteza de nul este că satisfacția cu viața este similară pentru toate cele trei grupuri. Totuși, noi ne așteptăm să apară diferențe : cei care au o relație ar trebui să aibă o satisfacție cu viața mai mare decât în cazul celorlalți. ANOVA este obținută din meniul **Analyze > Compare Means > One-Way ANOVA** (figura 6.8). Fereastra este la fel de intuitivă ca la testul t pentru eșantioane independente. În partea stângă, avem lista de variabile din care le vom selecta pe cele utilizate în analiză. La **Dependent List** introducem variabila metrică, pentru care calculăm mediile. Aici se introduce satisfacția cu viața care poartă numele V23. La **Factor** introducem variabila categorială, cea care dă grupurile pe care dorim să le comparăm după nivelul satisfacției cu viața. Aici se introduce starea civilă care poartă numele v57rec. Spre deosebire de meniul testului t pentru eșantioane independente, aici nu mai este nevoie să definim codurile grupurilor. Pentru a fi relevantă comparația, fiecare grup trebuie să aibă un număr decent de cazuri. Dacă nu se întâmplă acest lucru, atunci este utilă combinarea lor. De exemplu, primele două categorii ale stării civile sunt obținute prin combinarea categoriilor din variabila inițială. Pe lângă asigurarea unui număr rezonabil de cazuri în fiecare grup, am avut în vedere și scopurile teoretice ale analizei.

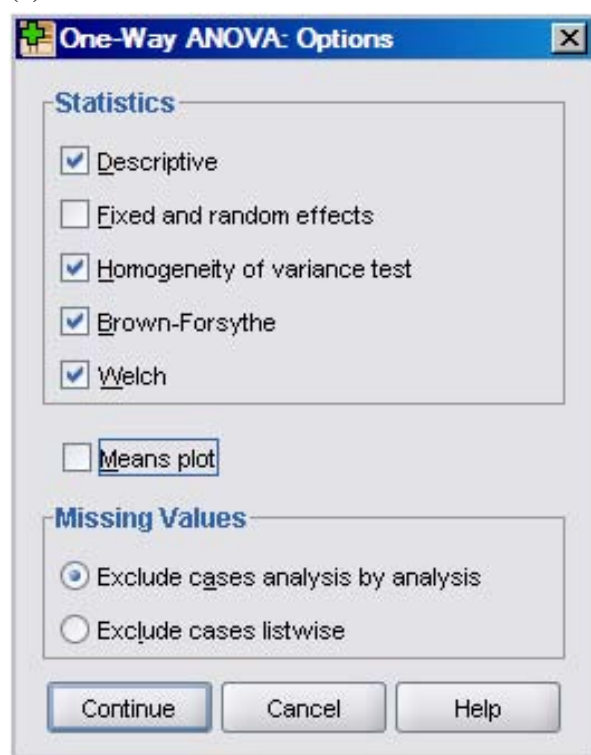
Dacă apăsăm butonul **Options**, putem alege mai multe opțiuni care vor fi afișate în **Output**. **Descriptive** ne oferă numărul de persoane din fiecare grup, media și abaterea standard a satisfacției cu viața, dar și intervalele de încredere în jurul mediilor și nu numai. **Homogeneity of variance test** ne oferă testul Levene. **Brown-Forsythe** și **Welch** sunt alternative robuste la testul F clasic specific ANOVA, atunci când asumția egalității varianțelor nu este îndeplinită. Dacă apăsăm **Continue** și **OK**, obținem rezultatul din tabelul 6.13a. Testul Levene ne spune că varianțele nu sunt egale : p calculat este mai mic decât 0.01 (coloana **Sig**). Acest lucru îl intuiau după ce am comparat abaterile standard. Putem să mai consultăm forma distribuției satisfacției cu viața pentru cele trei grupuri realizând câte un grafic bară pentru fiecare dintre ele. Dacă pentru cei care au o relație și cei care sunt divorțați, separați sau singuri distribuțiile au aproximativ aceeași formă, aceasta arată destul de diferit pentru văduvi. Testul **F**, din tabelul **ANOVA**, ne spune că cel puțin două dintre grupurile comparate au o satisfacție cu viața diferită. Pentru că varianțele sunt inegale, am preferat să consult și alternativele **Welch** și **Brown-Forsythe** care, de data aceasta, oferă același rezultat ca și testul **F**. Pentru a afla care grupuri diferă și în ce fel, trebuie să folosim metoda comparațiilor multiple sau testele post-hoc. Acestea pot fi accesate apăsând butonul **Post Hoc**. Aceste teste, fiecare cu avantajele și dezavantajele sale, sunt alese în funcție de rezultatul testului Levene. Aici, pentru că varianțele nu sunt egale, alegem unul dintre testele din secțiunea **Equal Variances Not Assumed** (figura 6.8c). Rezultatele sunt prezentate în tabelul 6.13b.

**Figura 6.8.** Meniul One-Way ANOVA. Analiza de varianță

(a)



(b)





(c)

**One-Way ANOVA: Post Hoc Multiple Comparisons**

**Equal Variances Assumed**

☐ LSD ☐ S-N-K ☐ Waller-Duncan  
☐ Bonferroni ☐ Tukey Type I/Type II Error Ratio: 100  
☐ Sidak ☐ Tukey's-b ☐ Dunnett  
☐ Scheffe ☐ Duncan Control Category: Last  
☐ R-E-G-WF ☐ Hochberg's GT2  
☐ R-E-G-WQ ☐ Gabriel

**Test**  
☒ 2-sided ☐ < Control ☐ > Control

**Equal Variances Not Assumed**

☒ Tamhane's T2 ☐ Dunnett's T3 ☐ Games-Howell ☐ Dunnett's C

Significance level: 0.05

Continue Cancel Help

Tabelul 6.13. Rezultate ale analizei de varianță

(a)

| Test of Homogeneity of Variances |     |      |      |
|----------------------------------|-----|------|------|
| V23 Satisfaction with your life  |     |      |      |
| Levene Statistic                 | df1 | df2  | Sig. |
| 13.805                           | 2   | 1482 | .000 |

| ANOVA                           |                |      |             |        |      |
|---------------------------------|----------------|------|-------------|--------|------|
| V23 Satisfaction with your life |                |      |             |        |      |
|                                 | Sum of Squares | df   | Mean Square | F      | Sig. |
| Between Groups                  | 284.112        | 2    | 142.056     | 25.670 | .000 |
| Within Groups                   | 8201.352       | 1482 | 5.534       |        |      |
| Total                           | 8485.464       | 1484 |             |        |      |

| Robust Tests of Equality of Means |                        |     |         |      |
|-----------------------------------|------------------------|-----|---------|------|
| V23 Satisfaction with your life   |                        |     |         |      |
|                                   | Statistic <sup>a</sup> | df1 | df2     | Sig. |
| Welch                             | 18.075                 | 2   | 365.203 | .000 |
| Brown-Forsythe                    | 21.310                 | 2   | 451.001 | .000 |

a. Asymptotically F distributed.



(b)

| <b>Multiple Comparisons</b>                     |   |                       |            |      |                         |             |
|---|---|-----------------------|------------|------|-------------------------|-------------|
| V23 Satisfaction with your life<br>Tamhane      |   |                       |            |      |                         |             |
| (I) v57rec starea civila (recodificare din V57) | (J) v57rec starea civila (recodificare din V57) | Mean Difference (I-J) | Std. Error | Sig. | 95% Confidence Interval |             |
|   |   |                       |            |      | Lower Bound             | Upper Bound |
| 1 casatorit sau are o relatie                   | 2 divortat, separat sau singur                  | .015                  | .146       | .999 | -.34                    | .37         |
|   | 3 vaduv   | 1.420*                | .238       | .000 | .85                     | 1.99        |
| 2 divortat, separat sau singur                  | 1 casatorit sau are o relatie                   | -.015                 | .146       | .999 | -.37                    | .34         |
|   | 3 vaduv   | 1.405*                | .260       | .000 | .78                     | 2.03        |
| 3 vaduv   | 1 casatorit sau are o relatie                   | -1.420*               | .238       | .000 | -1.99                   | -.85        |
|   | 2 divortat, separat sau singur                  | -1.405*               | .260       | .000 | -2.03                   | -.78        |

\*. The mean difference is significant at the 0.05 level.

Mai întâi, consultăm coloana **Sig** care conține valorile  $p$  calculate. Observăm diferențe semnificative statistic ( $p$  calculat  $< 0.05$ ) între grupurile „căsătorit sau are o relație” și „văduv”, respectiv „divorțat, separat sau singur” și „văduv”. Ipoteza de lucru se confirmă parțial, pentru că nu observăm o diferență între „căsătorit sau are o relație” și „divorțat, separat sau singur”. Apoi consultăm coloana **Mean Difference (I-J)**. Aceasta ne spune cu cât diferă mediile grupurilor comparate. Literele **I** și **J** desemnează prima, respectiv a doua coloană din tabel. De exemplu, diferența dintre media satisfacției cu viața a celor căsătoriți sau care au o relație și media văduvilor este de 1.42 unități. Cei dintâi au media 6.85, iar cei din urmă au media 5.43. Pasul următor este căutarea acelor factori care fac, de exemplu, ca persoanele divorțate, separate sau singure să fie mai satisfăcute cu viața decât cele văduve.

## 6.4. Două grafice uzuale în descrierea datelor

Graficele pot fi folosite în două scopuri. Un scop, pur operațional, este vizualizarea datelor altfel decât sub formă de tabele în timpul activității de analiză. Celălalt scop este prezentarea informației în lucrările noastre într-un mod mai intuitiv decât sub formă de tabele sau descriere în cuvinte. A face un grafic bun nu este atât de intuitiv tot timpul.

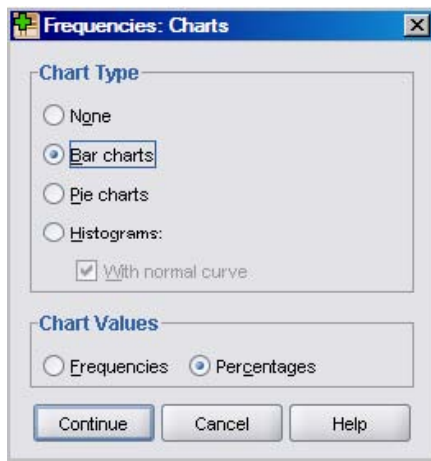
Good și Hardin (2012) oferă câteva reguli pentru cei ce doresc să utilizeze grafice pentru prezentarea datelor. Prezentăm o listă adaptată după acești autori :

- Realizați grafice 2D. A treia dimensiune trebuie folosită doar dacă există.
- Folosiți bare și evitați formele geometrice speciale cum ar fi conul, cilindrul etc. De asemenea, evitați umbrele generate de bare.
- Includeți valorile pe bare. Dacă sunt prea multe bare, includeți valorile măcar la extreme și la o categorie de interes major. În cazul în care comparăm salariul minim din diferite țări europene, evidențiem țara cu salariul minim și pe cea cu salariul maxim, dar și valoarea specifică României.
- Includeți etichete care clarifică elementele din grafic. Dacă sunt prea multe etichete, atunci alegeți cu atenție unele care evidențiază ideea centrală a graficului. Nu suprapuneți etichetele cu elementele esențiale ale graficului.
- Evitați spațiile goale de dimensiuni mari în grafice. Ajustați scala variabilelor astfel încât să reflecte amplitudinea din date nu pe cea ideală.
- Utilizați graficele în acord cu proprietățile variabilelor.

Această listă este doar un început. Poate fi îmbogățită și adaptată în funcție de informația care se vrea transmisă și graficul ales în acest sens.

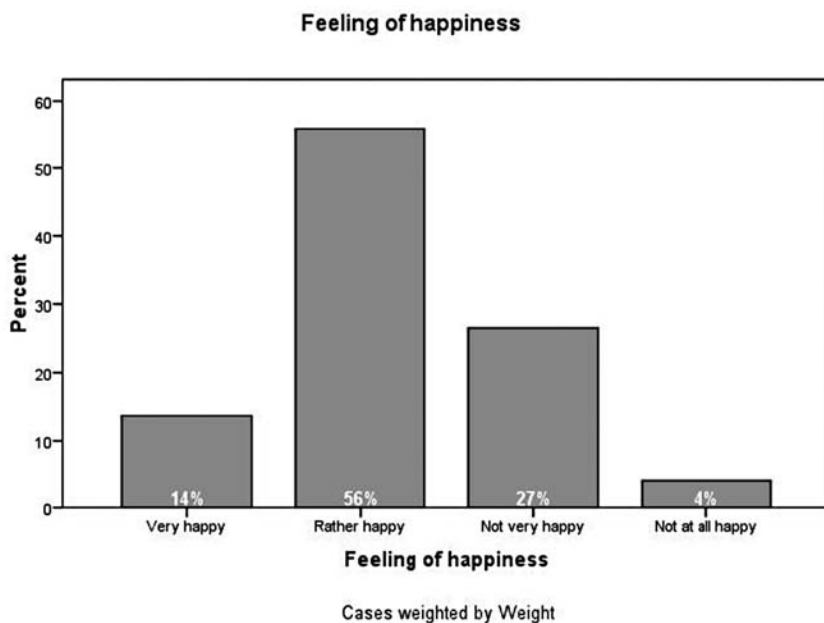
SPSS are o caracteristică utilă mai ales pentru utilizatorii novici : ne permite să realizăm grafice atât din meniurile unor analize, cât și din meniul dedicat special acestui lucru. Pentru începători, recomand utilizarea primei variante. De exemplu, meniul **Frequencies** ne oferă posibilitatea realizării a trei grafice : radial (**pie**), bară (**bar chart**) și histogramă (**histogram**). Intrați în meniu și apăsați butonul **Charts**. În fereastra care se deschide (figura 6.9) trebuie doar să selectăm tipul de grafic care ne interesează. Aici am ales să facem un grafic bară, axa Oy evidențiind procente. Folosim procente pentru că acestea au mai mult sens atunci când le citim decât frecvențele absolute.

Figura 6.9. Meniul Frequencies, Charts



Graficul radial nu este un instrument bun de vizualizare a datelor. Diferențele dintre secțiunile graficului pot fi atât de mici, încât devine necesară suprapunerea valorilor peste ele. Deja oferim informație redundantă. O soluție mult mai bună este graficul bară. Un exemplu este prezentat în figura 6.10. Pe axa Ox sunt cele patru niveluri de fericire. Axa Oy ne spune procentul celor care aleg o categorie sau alta.

**Figura 6.10.** Grafic bară



Valorile de pe fiecare bară le-am adăugat ulterior. Am dat dublu click pe grafic, acțiune în urma căreia se deschide **Chart Editor**. Mergem în meniul **Elements > Show Data Labels**. În fereastra care se deschide, selectăm tabul **Data Value Labels**, **Label Position** și apoi **Custom**, bifând poziția dorită (figura 6.11a).

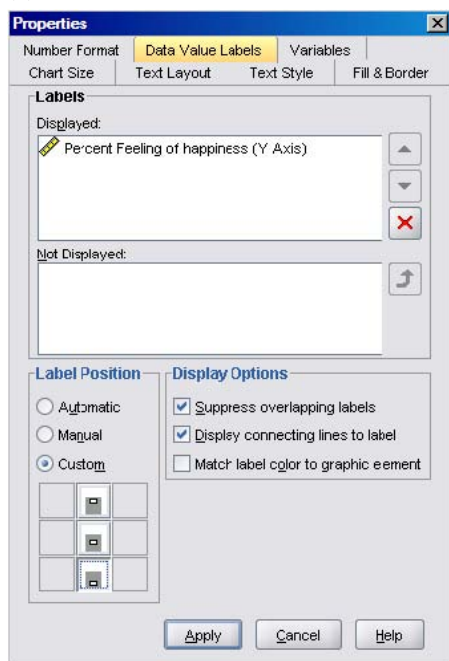
Tot aici putem modifica numărul de zecimale, tipul textului și altele. Selectăm tabul **Number Format** și facem următoarele modificări : **Decimal Places = 0** și **Trailing Characters = %** (figura 6.11b). Apăsăm **Apply**, **Close** și închidem fereastra de editare a graficului. Evident, putem face și alte modificări în aceste ferestre.

Aici, vedem rapid că majoritatea românilor erau destul de fericiți în 2012 conform WVS.

Când vrem să reprezentăm grafic o variabilă metrică cu multe categorii, în locul graficului bară alegem histograma (figura 6.12). Mergem în **Frequencies**, apăsăm **Charts**, și selectăm **Histogram**. Comparând cele două histograme, observăm o variație mai mare în rândul femeilor în ceea ce privește timpul petrecut cu îngrijirea copiilor, bătrânilor sau celor bolnavi. Graficele sugerează, de asemenea, un timp mai mare petrecut de femei cu acest gen de activități.

**Figura 6.11.** Editarea valorilor de pe bare, Properties

(a)



(b)

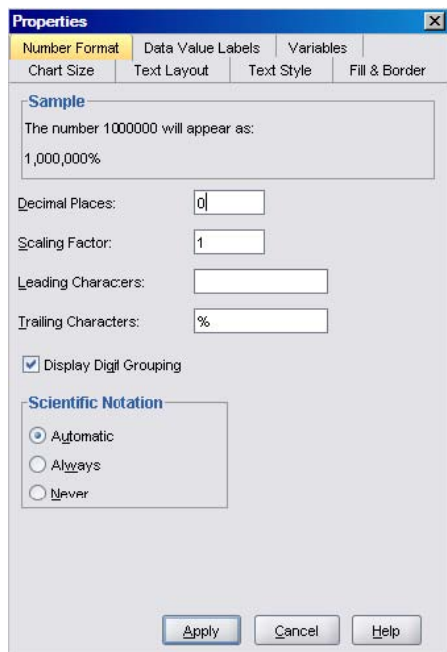
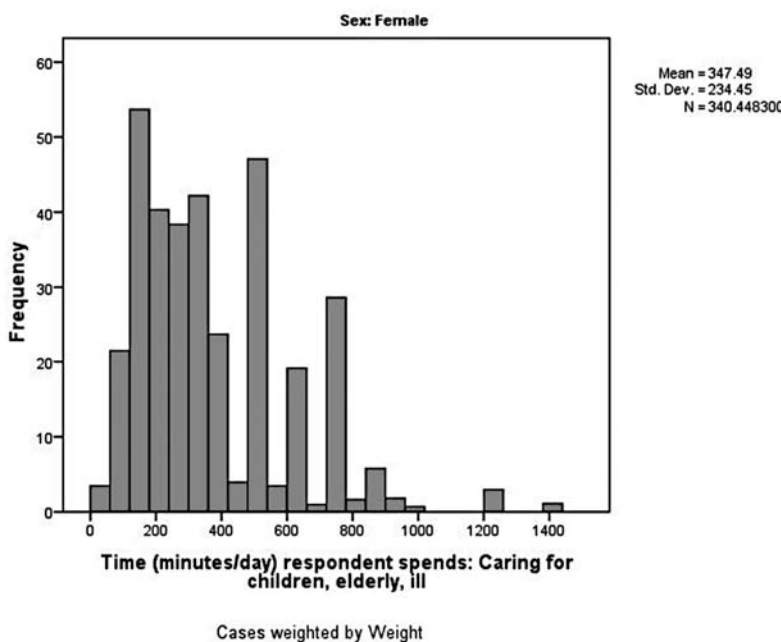
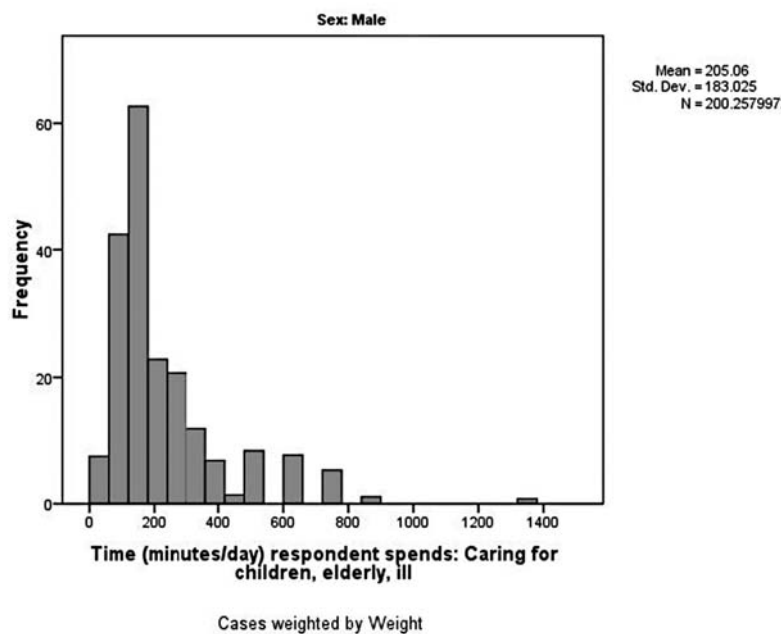
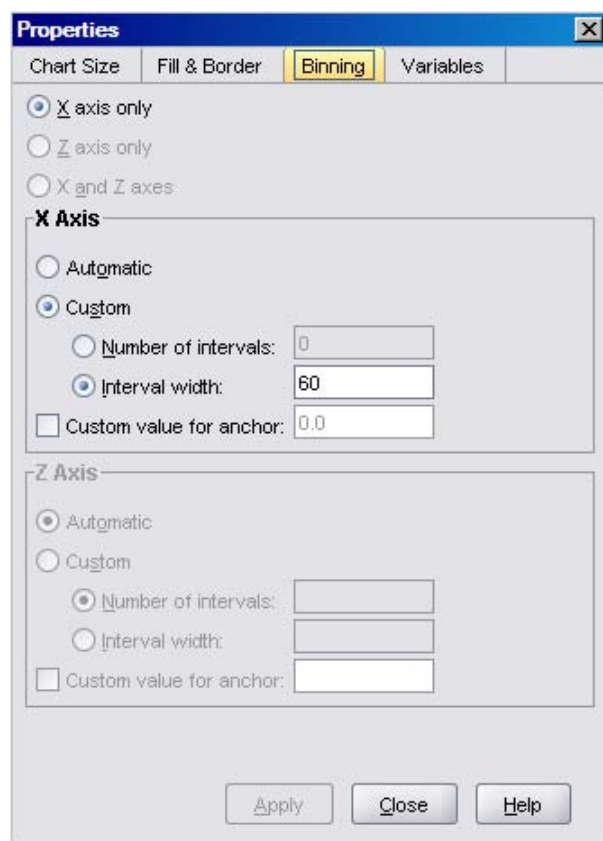


Figura 6.12. Histograme



Histogramele trebuie utilizate cu precauție. De exemplu, aici, am modificat mărimea intervalului folosit pentru reprezentarea grafică astfel încât să fie egal cu 60 de minute. Acest lucru se poate face dând, în **Output**, dublu click pe grafic. Se deschide **Chart Editor**. În interiorul acestuia, dăm dublu click pe barele histogramei și se deschide fereastra **Properties**. Selectăm tabul **Binning** (figura 6.13) și, în secțiunea **X Axis**, selectăm **Custom**, iar la **Interval width** introducem valoarea dorită.

**Figura 6.13.** Chart Editor, Properties pentru histogramă



Există mai multe lucrări care prezintă principiile reprezentării grafice corecte. Dintre acestea le pot aminti pe cele elaborate de Chambers și colaboratorii (1983), Jacoby (1997, 1998) și Tufte (2001).

## 6.5. Exerciții

Pentru aceste exerciții, utilizăm baza de date și/sau chestionarul World Values Survey 2012 rezultat(ă) în urma aplicării chestionarului în România. Baza de date poate fi descărcată de pe pagina de internet a *Grupului Românesc pentru Studiul Valorilor Sociale* (<http://www.romanianvalues.ro>).

1. Deschideți chestionarul WVS 2012. Calculați indicatorii tendinței centrale și variației corespunzători pentru variabilele de pe paginile cu număr impar.
2. Exportați tabelele în Excel și editați-le pentru includerea într-un raport.
3. Pentru variabilele de la exercițiul anterior, realizați câte un grafic care să reflecte cât mai bine informația. Editați aceste grafice astfel încât să poată fi folosite într-un material tipărit monocrom (alb-negru).
4. Găsiți în baza de date variabila care măsoară fericirea (nu satisfacția cu viața). Elaborați o listă cu zece variabile nominale, diferite față de cele folosite în textul acestui capitol, care credeți că influențează fericirea. Realizați o listă de ipoteze în care fericirea este variabila explicată.
5. Realizați zece tabele de contingență în care testați ipotezele notate la exercițiul anterior. Scrieți un scurt raport de o pagină în care descrieți ce ați aflat, folosind valorile reziduale ajustate.
6. Realizați un profil al părinților care sunt predispuși să transmită copiilor lor valoarea „imaginație”. Profilul trebuie să conțină cinci variabile explicative nominale sau ordinale din chestionarul WVS 2012.
7. Testați profilul folosind tabele de contingență cu valori reziduale ajustate.
8. Creați o variabilă nouă care reflectă statutul de membru activ în organizații voluntare. Verificați dacă media satisfacției cu viața este diferită pentru membrii activi față de cei care sunt membri inactivi sau nu sunt membri.
9. Creați o variabilă nouă care reflectă statutul de membru activ, membru inactiv și nonmembru în organizații voluntare. Verificați dacă media satisfacției cu viața este diferită pentru aceste trei categorii.
10. Creați o variabilă nouă care să reflecte intoleranța față de grupuri marginale sau minorități (persoane dependente de droguri, persoane de rasă diferită de a dumneavoastră etc.). Aceasta trebuie să folosească setul de variabile v36-v44 din chestionarul WVS 2012. Variabila nou-creată trebuie să reprezinte numărul de grupuri pe care o persoană nu i-ar dori ca vecini. Verificați dacă numărul este mai mic pentru cei care fac voluntariat activ decât pentru cei care nu fac voluntariat. Verificați dacă există diferențe în ceea ce privește numărul de categorii alese între cei care fac voluntariat activ, cei care nu fac voluntariat activ, respectiv cei care nu fac voluntariat deloc.





## 7. Explorarea datelor : asumptii

Dacă două persoane au aceleași caracteristici, dar diferă în funcție de venitul lunar câștigat, care va fi mai satisfăcută cu viața : cea cu un venit mai mare ? În orașele mari sunt mai multe persoane fericite decât în orașele mici ? Persoanele care au emigrat la o vârstă mai înaintată se adaptează mai ușor la modul de viață din țara de destinație ? Studenții care participă la activități de voluntariat au o șansă mai mare în a găsi un loc de muncă apropiat de așteptările și dorințele lor ?

Un cercetător organizat, înainte de a trece la elaborarea chestionarului, trebuie să anticipeze ce fel de analize solicită întrebările sale de cercetare. Deși pare contraintuitiv, punând căruța înaintea calului, în practică nu este chiar așa. Degeaba formulezi întrebările și alegi variantele de răspuns dacă nu ești conștient de calitățile psihometrice pe care acestea le au. Când primește datele din teren și ai în față baza de date, începi să te întrebi ce poți face cu variabilele avute la dispoziție și, dacă nu te-ai gândit dinainte la corespondența cu tehnicile statistice, răspunsul s-ar putea să nu îți placă. Să presupunem că, măcar parțial, cercetătorul a formulat întrebările din chestionar în acord cu cerințele statistice ale obiectivelor de cercetare. Înainte de a trece la analiza propriu-zisă a datelor, adică la aplicarea tehnicii statistice prin care acestea răspund la întrebarea de cercetare, va trebui să inspecteze variabilele univariat, bivariat sau chiar multivariat. Când explorăm datele, căutăm să înțelegem dacă variabilele au suficientă variație, dacă există cazuri extreme, cum arată distribuția acestora etc.

Pentru înțelegerea unor analize descrise aici, cititorul ar trebui să aibă cunoștințe minime despre testele de semnificație. Câteva noțiuni elementare din această zonă a statisticii au fost prezentate în capitolul 6. Totuși, ar trebui să suplimentați această lucrare cu un manual de statistică. Am sugerat câteva astfel de lucrări în volumul de față.

În acest capitol, vom vorbi despre statistici și grafice care ne ajută să decidem dacă există cazuri extreme, care este forma distribuției, dacă între două variabile există o relație liniară etc. Acestea ne ajută să ne cunoaștem datele înainte de a le folosi la calcularea unor statistici cum ar fi media sau abaterea standard. Analistii trebuie să fie sceptici în legătură cu calitatea datelor lor. Acest scepticism previne erorile în interpretarea substanțială.

## 7.1. Distribuția unei variabile

Măsurarea este procesul prin care cercetătorul operaționalizează conceptele cu care lucrează, construiește instrumentele prin care culege date despre fenomenele pe care le cuprinde și verifică validitatea și fidelitatea acestora. Prin măsurare, de exemplu, ajungem la scalele compuse pe care le folosim frecvent în chestionarele noastre. Unul dintre criteriile după care cercetătorul decide care sunt itemii pe care îi va păstra în scala finală este cel al irelevanței (Mărginean, 1982). Pe scurt, acest criteriu presupune eliminarea, din analizele prin care căutăm răspunsurile la întrebarea de cercetare, a itemilor cu care toți respondenții sau o majoritate covârșitoare dintre aceștia sunt de acord sau nu. Pentru variabilele categoriale este relativ simplu să aflăm această informație. Putem realiza un tabel de frecvență și inspecta distribuția procentelor (tabelul 7.1).

**Tabelul 7.1.** Tabel de frecvență : verificarea variației variabilelor categoriale

| V219 Information source: TV news |                     |           |         |               |                    |
|----------------------------------|---------------------|-----------|---------|---------------|--------------------|
|                                  |                     | Frequency | Percent | Valid Percent | Cumulative Percent |
| Valid                            | 1 Daily             | 1235      | 82.1    | 82.3          | 82.3               |
|                                  | 2 Weekly            | 137       | 9.1     | 9.1           | 91.4               |
|                                  | 3 Monthly           | 19        | 1.3     | 1.3           | 92.6               |
|                                  | 4 Less than monthly | 81        | 5.4     | 5.4           | 98.0               |
|                                  | 5 Never             | 30        | 2.0     | 2.0           | 100.0              |
|                                  | Total               | 1501      | 99.9    | 100.0         |                    |
| Missing                          | -2 No answer        | 1         | .0      |               |                    |
|                                  | -1 Don't know       | 1         | .1      |               |                    |
|                                  | Total               | 2         | .1      |               |                    |
| Total                            |                     | 1503      | 100.0   |               |                    |

Conform datelor WVS 2012, 82% dintre români foloseau, zilnic, știrile prezentate la televizor ca sursă de informare despre ce se petrece în țară și în lume. Aceasta este o informație utilă și interesantă despre comportamentul de informare al românilor și nu numai. Are însă suficientă variație această variabilă dacă dorim să o includem într-o analiză multivariată? Răspunsul nu este simplu de oferit. Dacă cercetătorul are argumente teoretice solide, poate decide să o folosească ca atare sau poate considera că ar fi mai util să o recodifice: televizorul este, probabil, principala sursă de informație și divertisment, fiind accesibil atât în ceea ce privește costurile, cât și referitor la dificultatea conținuturilor prezentate. Presiunea timpului solicită formate scurte, concentrate, cu mesaje transmise în forme inteligibile pentru mase mari de privitori. Din acest

motiv, orice frecvență de utilizare, în afara celei zilnice, ar putea fi considerată aparte. Cei care se informează mai rar decât zilnic de la știrile televizate au, probabil, caracteristici diferite față de ceilalți. Viața socială este atât de complexă, încât decizia de a utiliza într-o formă sau alta această variabilă depinde de mulți factori.

Să nu înțelegeți că procesul de explorare a datelor are ca unic scop găsirea problemelor. În primul rând, dorim să ne familiarizăm cu datele. Apoi, dorim să vedem dacă sunt modificări pe care trebuie să le aducem variabilelor pentru a utiliza cât mai multă informație culeasă prin chestionar.

Agresti și Finlay (2008) consideră distribuția normală (cea care are formă de clopot) ca fiind cea mai importantă pentru analiza statistică deoarece aproximează destul de bine forma multor variabile din viața reală. Pentru a înțelege proprietățile acestei distribuții și de ce este importantă, trebuie să consultați capitolele dedicate acestui subiect din lucrarea citată sau din Agresti și Franklin (2013). Deși, în esență, cele două lucrări prezintă aceeași informație, cea din urmă are o formă de prezentare grafică mai prietenoasă. În acest moment, este suficient să reținem că multe dintre analizele inferențiale aplicate în mod obișnuit în științele sociale folosesc această distribuție. Concepte importante asociate sunt distribuția de eșantionare, eroarea standard sau teorema limită centrală. Toate sunt tratate comprehensiv în lucrările citate.

Testele parametrice asumă distribuția normală. De aceea, înainte de a rula și interpreta astfel de teste, trebuie să verificăm dacă variabilele au sau nu o distribuție aproximativ normală. Dacă nu au, atunci putem utiliza testele nonparametrice echivalente (Cramer și Howitt, 2004). Printre analizele pentru care verificarea asumției normalității distribuției este importantă se numără analiza de corelație Pearson, ANOVA, testele *t* și regresia multiplă (de Vaus, 2002). Această discuție se aplică variabilelor metrice. În practică, convențional, se acceptă și pentru variabilele ordinale (scală tip Likert).

Putem verifica dacă această asumție este îndeplinită folosind grafice sau calculând diferiți indicatori și teste statistice. Mai întâi discutăm despre metodele grafice.

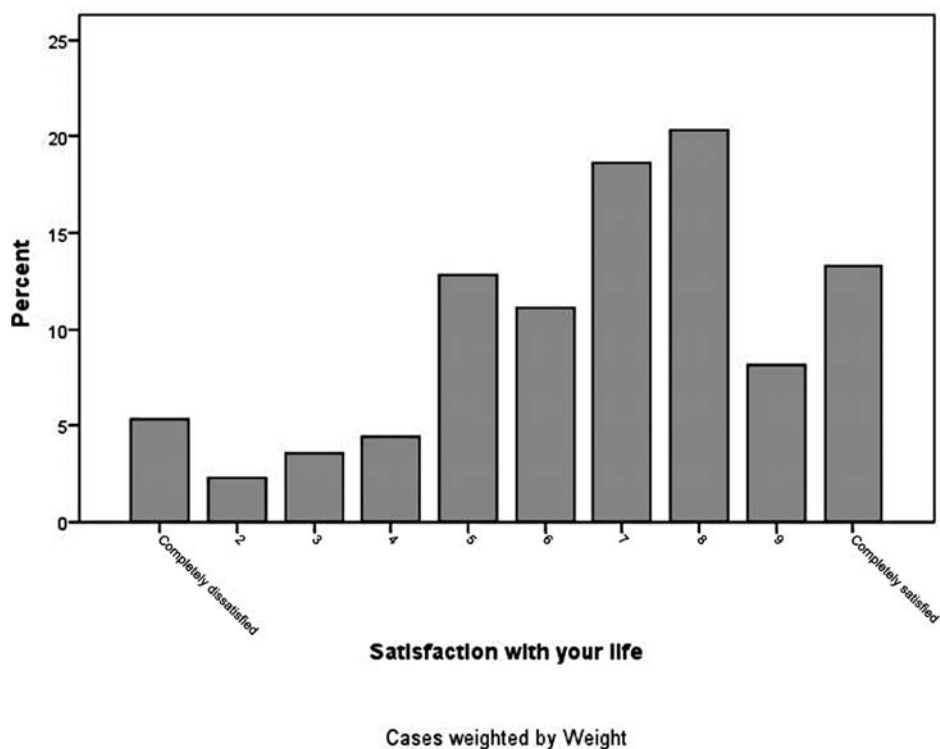
În figura 7.1 este prezentată distribuția satisfacției cu viața pentru români în 2012 conform World Values Survey.

Acest grafic a fost obținut din meniul **Analyze > Descriptive statistics > Frequencies > Charts**. A fost editat aplicând pașii deja discutați în alte locuri din acest volum.

Distribuția satisfacției cu viața, o variabilă măsurată pe o scală de la 1 la 10, unde 1 înseamnă „total nemulțumit” și 10 „total mulțumit”, este alungită la stânga. Majoritatea românilor se poziționează în partea pozitivă a scalei. Distribuția se abate de la normalitate. Acest lucru nu este neapărat rău, pentru că satisfacția cu viața, în realitate, are o distribuție de acest gen (Cummins, 2003). Ceea ce

îngrijorează, gândindu-ne la validitatea rezultatelor analizelor statistice, sunt frecvențele neașteptat de mari pentru categoriile 1 și 10, cei total nemulțumiți sau mulțumiți cu viața lor în general. Înainte de a calcula medii sau coeficienți de corelație, trebuie să înțelegem de ce apar aceste două abateri de la normalitate. O altă problemă importantă, care este mai puțin vizibilă aici, apare atunci când există mai multe vârfuri și goluri între aceste vârfuri. Distribuția, în această situație, apare ca și când ar fi formată din mai multe distribuții mici. Calcularea mediei sau medianei nu ar avea foarte mult sens în această situație. Hartwig și Dearing (1979) recomandă să creăm variabile separate din cea inițială și să le utilizăm ca atare.

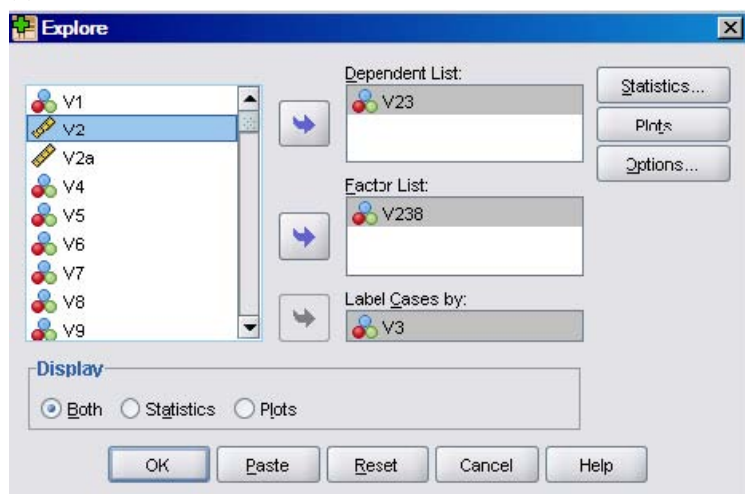
**Figura 7.1.** Grafic bară pentru verificarea asumției de normalitate a distribuției



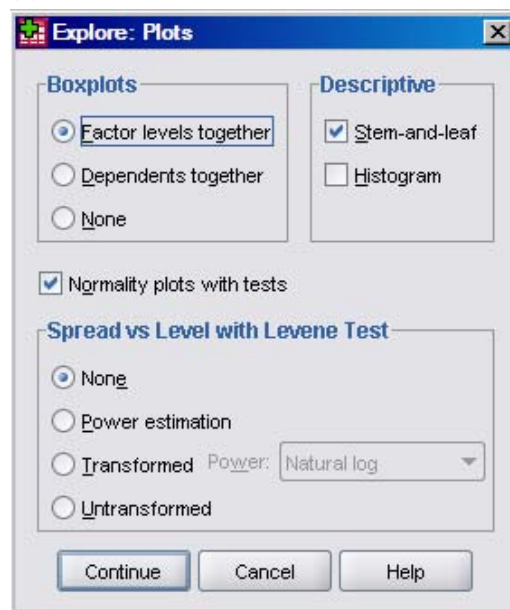
Un grafic care oferă informații similare cu histograma, dar care adaugă și altele noi, este box-plot-ul. Acesta poate fi obținut din **Analyze > Descriptive statistics > Explore > Plots > Box-plot = Factor levels together** (figura 7.2).

Figura 7.2. Meniul Explore

(a)

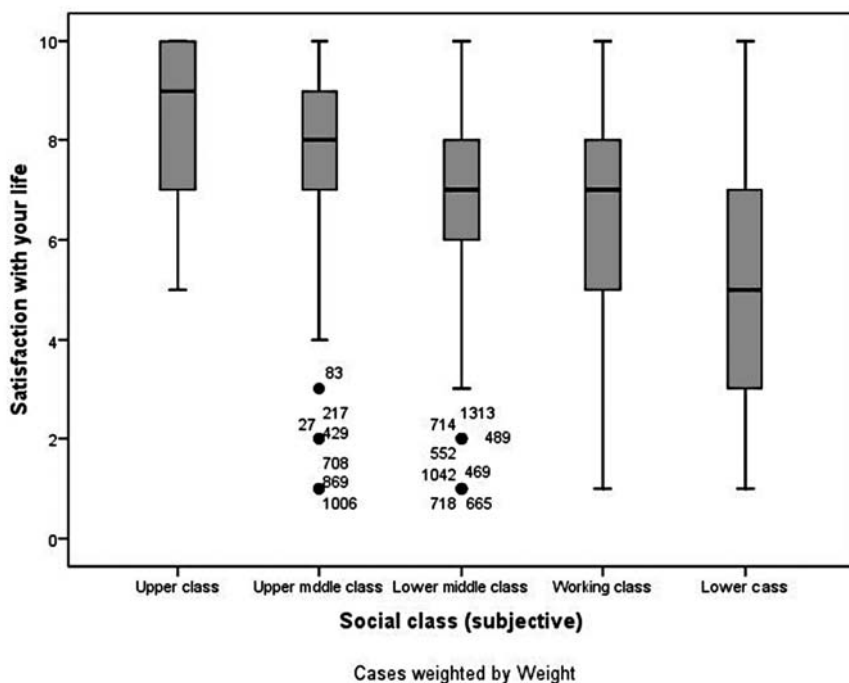


(b)



În WVS 2012, respondenții au fost rugați să se poziționeze în ierarhia socială într-una dintre pozițiile : „clasa de jos”, „clasa muncitoare”, „în partea de jos a clasei mijlocii”, „în partea de sus a clasei mijlocii”, „în clasa de sus”. În figura 7.3 este prezentată distribuția satisfacției cu viața pentru fiecare dintre aceste poziții sociale.

**Figura 7.3.** Box-plot : distribuția satisfacției cu viața în funcție de poziția socială subiectivă




Revenind puțin la meniul din care a fost obținut box-plot-ul, observăm următoarele (figura 7.2a) :

- Presupunem că o variabilă, aici satisfacția cu viața (V23), variază în funcție de o alta, aici autopozitionarea pe o scală a poziției sociale (V238). De aceea, variabila pentru care sunt calculate statisticile (medie, mediană, abatere standard etc.) va fi introdusă în câmpul **Dependent List**, iar variabila care dă grupurile pentru care sunt realizate comparațiile va fi introdusă în câmpul **Factor List**. Denumirea **factor** atribuită unei variabile categoricale care distinge între anumite grupuri a mai fost întâlnită la analiza de varianță.
- Pentru a putea identifica mai ușor eventualele cazuri cu „probleme”, le etichetăm folosind variabila de identificare care nu trebuie să lipsească din nici o bază de date. Aici, această variabilă, care conține ID-uri unice pentru fiecare respondent, este **V3**.
- Pentru a nu încărca outputul cu multe informații, în prima fază, putem bifa doar opțiunea **Plots** în secțiunea **Display**. Dacă lăsăm bifată opțiunea **Both**, atunci în output vor fi afișate și statisticile produse de meniu. Aici acestea vor fi calculate pentru fiecare poziție socială. Am putea lua în considerare suprimarea temporară a acestora.

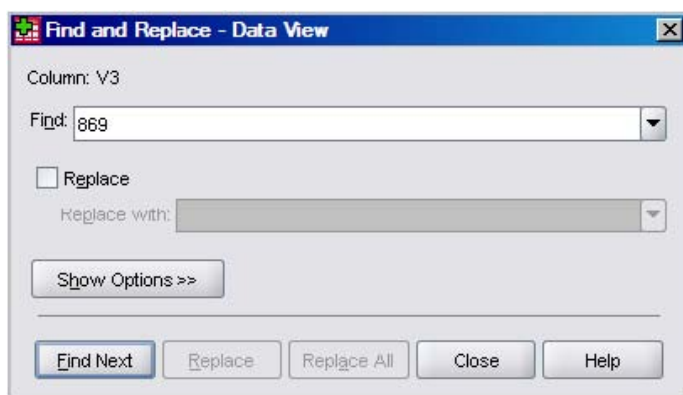
Spuneam că box-plot-ul ne oferă câteva informații mai puțin evidente din histogramă. Comparând pozițiile liniei îngroșate din interiorul cutiilor (mediana), observăm că aceasta ia valori din ce în ce mai mici pe măsură ce persoanele se autopoziționează în clase sociale aflate pe trepte din ce în ce mai joase ale ierarhiei sociale. În cazul de față trebuie să fim precauți cu interpretarea valorilor pentru „clasa de sus”, deoarece în această categorie sunt doar 17 persoane în eșantion. Apoi observăm că, mergând de la clasa de sus în jos, crește și variația satisfacției cu viața (cutia este mai lungă, deci abaterea intercuartilă este mai mare). Precauția cu privire la numărul de cazuri din clasa de sus capătă și mai mult sens aici, pentru că ne așteptam să vedem variația cea mai restrânsă în această categorie. Continuăm cu citirea graficului. În clasa muncitoare, distribuția este alungită la stânga (mediana se află înspre capătul de sus, cuartila 3, a cutiei). În cele două secțiuni ale clasei mijlocii, există cazuri extreme (outlieri) simbolizate prin cer-culețe. SPSS identifică două tipuri de cazuri extreme : cele discutate și cele care se află foarte departe în distribuție, reprezentate cu stelute. Atunci când încercăm să remediem problema cazurilor extreme, întotdeauna începem cu stelutele.

Tabachnick și Fidell (2007) oferă mai multe soluții pentru gestionarea cazurilor extreme. Prima soluție, cea radicală, este scoaterea din analize a persoanei sau a persoanelor respective. A doua soluție, cea care caută să maximizeze utilizarea datelor aflate la dispoziție, presupune aplicarea unor transformări variabilei care conține cazurile extreme. O astfel de transformare poate fi obținută, de exemplu, prin logaritmarea variabilei cu cazuri extreme. Cazurile vor rămâne în baza de date, dar influența lor va fi diminuată considerabil. O altă metodă constă în intervenții directe asupra cazurilor extreme : valoarea extremă este recodificată în jos sau în sus. De exemplu, dacă salariul din ultima lună are valoarea extremă de 15.000 de lei, iar următoarea valoare, care nu este caz extrem, este 5.500 de lei, atunci putem recodifica în 5.600 de lei sau altă valoare aleasă în funcție de distribuția celorlalte valori. Soluțiile nu sunt simplu de ales. Lucrurile se complică și mai mult dacă reținem că discuția, până în acest punct, a fost despre cazurile extreme univariate ignorându-le pe cele multivariate (o persoană care a absolvit facultatea, are 24 de ani și la primul loc de muncă primește un salariu lunar de 10.000 de lei). Aceiași autori atrag atenția că soluțiile enunțate s-ar putea să nu funcționeze bine în modelele multivariate. Mai mult, Hair și colaboratorii (2010) atrag atenția că ștergerea sau modificarea cazurilor extreme poate avea un efect pervers grav : modelele multivariate vor fi mai bune din punct de vedere statistic, dar mai puțin generalizabile la populația pentru care facem inferențe. Vă recomand să consultați aceste două lucrări, pentru că oferă exemple detaliate despre cum se identifică și gestionează cazurile extreme.

Deoarece depășește scopul acestei lucrări, nu mai insistăm asupra aspectelor teoretice ale acestei teme, așadar vom discuta în continuare doar despre partea operațională a identificării rapide a cazurilor extreme. Am discutat despre cazurile extreme univariate. În figura 7.3 am observat că o persoană care are id-ul 869 este considerată caz extrem. În SPSS, putem consulta imediat răspunsurile oferite de acest respondent la diferite variabile din chestionar. O variantă presupune să

mergem în **Data View**, unde dăm click în prima celulă din dreptul variabilei V3 (cea care conține id-urile respondenților), apăsăm pe iconița  și tastăm în câmpul **Find** 869 (figura 7.4). Apoi apăsăm butonul **Find Next**, comandă care ne va duce la celula din V3 care conține numărul 869.

**Figura 7.4.** Find



Putem selecta rândul care conține acest respondent dând click pe numărul rândului, după care vom naviga folosind bara orizontală de scroll. Mai simplu ar fi însă ca ID-ul să rămână vizibil, schimbând doar poziția celorlalte variabile. Pentru aceasta, mergem în meniul **Window > Split** (figura 7.5).

**Figura 7.5.** Meniul Window > Split

| V3  | V4 | V6 | V7 | V8 | V9 |
|-----|----|----|----|----|----|
| 869 | 1  | 1  | 4  | 3  | 2  |
| C70 | 1  | 2  | 3  | 2  | 3  |
| E71 | 1  | 1  | 3  | 1  | 1  |
| E72 | 1  | 3  | 2  | 2  | 3  |

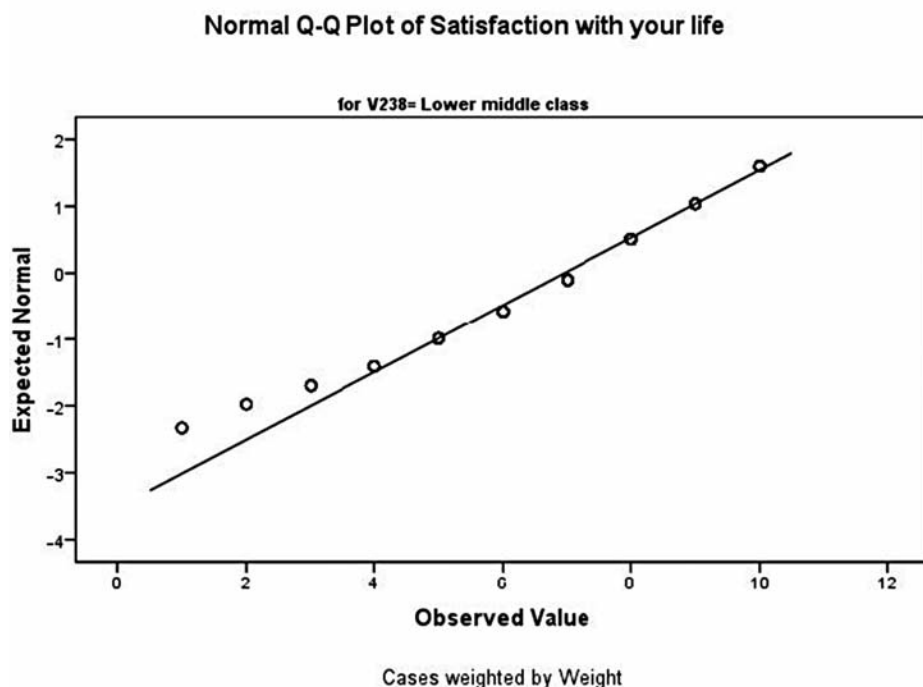
Fereastra se împarte în două sau patru secțiuni, în funcție de varianta de SPSS cu care lucrați, fiecare fiind navigabilă de sine stătător. Acum, de exemplu, ținând constantă poziția lui V3, în cadranul din stânga sus, putem naviga orizontal în cadranul din dreapta pentru a vedea ce valori ia respondentul cu ID-ul 869 la alte variabile. Există posibilitatea ca acea valoare extremă să fie doar o eroare de introducere a datelor. Adică operatorul de introducere, în loc să tasteze valoarea 5, a tastat valoarea 1. La variabilele subiective (valori, atitudini, evaluări etc.) este greu să ne dăm seama de aceste lucruri, dar la venituri sau proprietăți s-ar putea să fie mai ușor. Cert este că, înainte de a lua o decizie de transformare sau ștergere, trebuie să verificăm chestionarele în original (dacă putem).



Am insistat asupra cazurilor extreme pentru că, de multe ori, asumptia de normalitate a distribuției este încălcată pentru că acestea există.

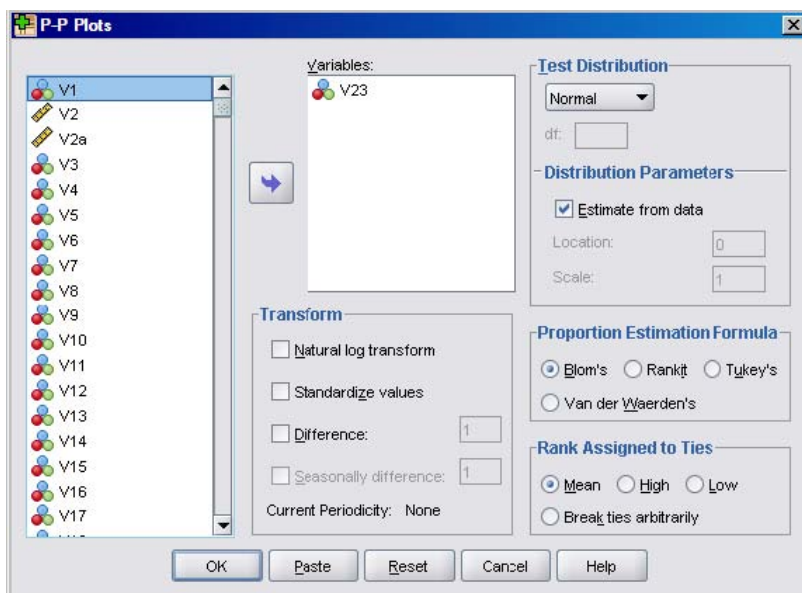
Histograma și box-plot-ul sunt rapid de construit și ușor de interpretat. Un alt grafic, creat special pentru evaluarea acestei asumptii, și care nu are dezavantajele celorlalte două, este **normal probability plot**. În figura 7.5 este prezentat un grafic similar, **normal q-q plot**, pentru satisfacția cu viața în cadrul categoriei „partea de jos a clasei mijlocii”. Distribuția normală este reprezentată prin linia diagonală, iar distribuția datelor din eșantion pentru satisfacția cu viața este dată de succesiunea cerculețelor. Dacă variabila are o distribuție normală, atunci cerculețele ar trebui să cadă aproximativ pe linie.

Figura 7.6. Normal probability plot

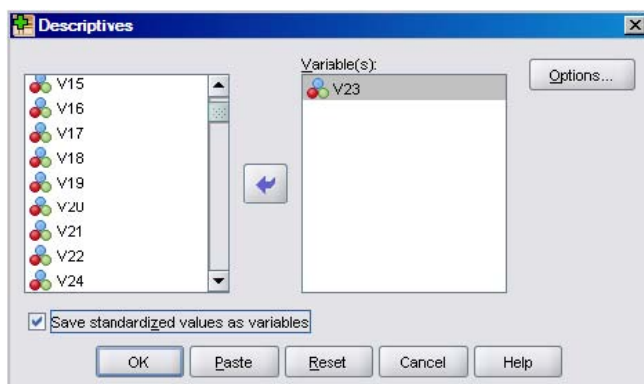


Acest grafic se obține din meniul **Analyze > Descriptive statistics > Explore > Plots = Normality plots with tests** (Figura 7.2b). Dacă dorim să vizualizăm graficele **normal p-p**, nu pe cele **normal q-q**, atunci putem folosi alt meniu pentru a le obține: **Analyze > Descriptive statistics > P-P Plots** sau **Q-Q Plots** (Figura 7.7). Dacă dorim să realizăm aceste grafice pentru diferite categorii ale altei variabile, așa cum am făcut cu autopозиționarea în ierarhia socială, atunci trebuie mai întâi să separăm (**split file**) baza de date după această variabilă sau să activăm anumite filtre (**select cases**).

Figura 7.7. Normal P-P Plots

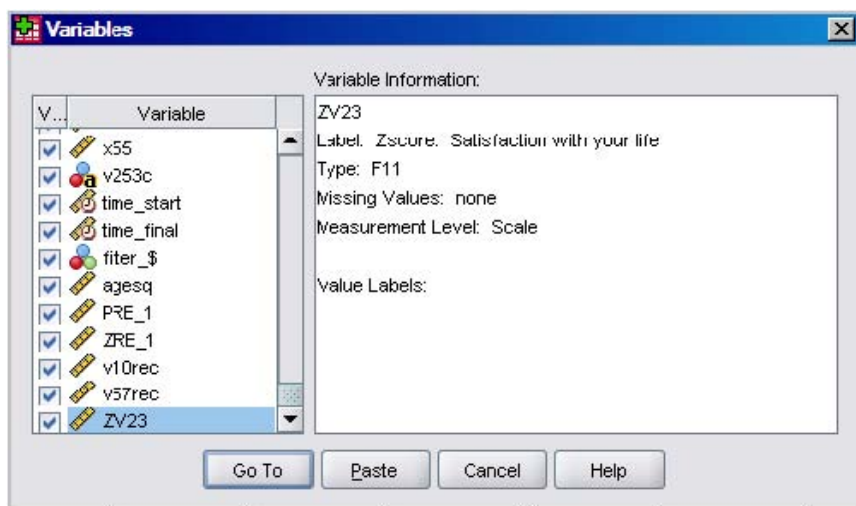


Graficele ne ajută, dar, deseori, avem nevoie și de indicatori specifici creați pentru aceleași scopuri. Aș vrea să ne întoarcem puțin la cazurile extreme. O metodă prin care verificăm dacă o valoare este extremă constă în transformarea acelei valori în scor  $z$ . Cazurile care au scoruri  $z$  la variabila explorată cu valori mai mari decât aproximativ 3, ignorând semnul, sunt potențiali outliers. Trebuie să ne reamintim aici de regula empirică aplicabilă distribuțiilor aproximativ normale (Agresti și Finlay, 2008). Scorurile  $z$  pot fi calculate din meniul **Analyze > Descriptive statistics > Descriptives > Save standardized values as variables** (figura 7.8). Opțiunea aceasta creează o variabilă nouă.

Figura 7.8. Descriptives, scoruri  $z$ 

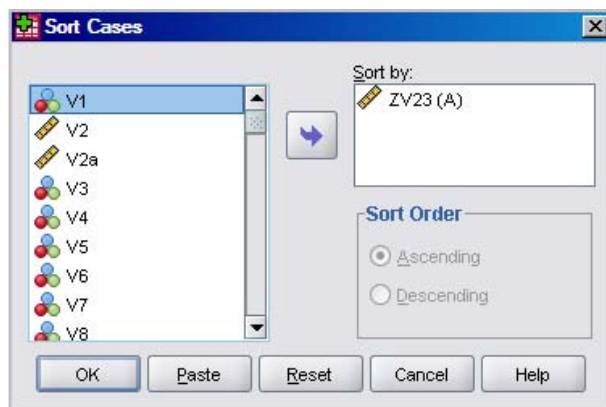
Pentru a vedea rezultatul, mergem în **Data View**, la variabila nou-creată. Aceasta va fi la sfârșitul bazei de date. Pentru a o găsi rapid, mergem în meniul **Utilities > Variables**. Se va deschide fereastra din figura 7.9.

Figura 7.9. Find variables



Selectăm orice variabilă din stânga și tastăm rapid litera z, deoarece variabila nou-creată va fi denumită automat de SPSS care va pune această literă ca prim caracter. Dacă avem mai multe variabile standardizate, atunci tastăm rapid zv23. Apoi apăsăm **Go To** și ne va duce la variabila dorită. Ne interesează valorile extreme, mai mari decât 3, în valoare absolută. Pentru a inspecta vizual ușor, sortăm baza de date. Mergem în meniul **Data > Sort Cases**, introducem variabila ZV23 în câmpul **Sort by** și, la **Sort Order**, bifăm **Ascending** (figura 7.10).

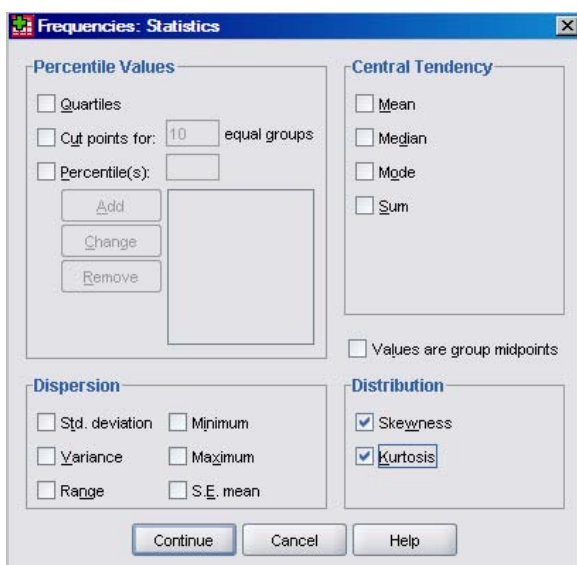
Figura 7.10. Sort Cases



Apoi mergem în **Data View** și inspectăm valorile. Cea mai mică este  $-2.37$ . Repetăm procedura, dar sortând **Descendind**. Cea mai mare este  $1.38$ . După criteriul  $z > 3$  (în valoare absolută), nu avem cazuri extreme. Observăm cum cercetătorul trebuie să își folosească rațiunea atunci când decide ce este și ce nu este caz extrem. Nu lăsăm programul să ia decizii pentru noi. Ca regulă de lucru, verificarea prin aplicarea mai multor tehnici cu același obiectiv este esențială.

Doi dintre cei mai utilizați indicatori statistici pentru verificarea normalității distribuției sunt **skewness** (alungire) și **kurtosis** (aplatizare). În SPSS, acești indicatori sunt centrați în jurul valorii 0 ce reprezintă distribuția normală. Când distribuția este alungită la dreapta, indicatorul de skewness are valoare pozitivă. Când distribuția este alungită la stânga, indicatorul de skewness are valoare negativă. Când observațiile sunt grupate strâns în jurul mediei, indicatorul de kurtosis are valoare pozitivă. Când observațiile sunt dispersate în jurul mediei, indicatorul de kurtosis are valoare negativă. Însă acești indicatori, în eșantioanele cu volum mare, pot să arate abateri de la normalitate, chiar și atunci când acestea sunt mici (Tabachnick și Fidell, 2007). De aceea, interpretarea lor trebuie combinată cu ceea ce ne oferă graficele discutate. Acești indicatori pot fi obținuți din mai multe meniuri: **Frequencies**, **Descriptives** sau **Explore**. Dacă ne interesează valorile pentru o variabilă în cadrul întregului eșantion, atunci putem să o alegem pe oricare dintre ele. Dacă vrem însă să aflăm aceste valori pentru diferite categorii ale altei variabile, ar fi mai util să folosim **Explore**. Motivul este simplu: pentru **Frequencies** și **Descriptives** ar trebui, în prealabil, să separăm (**split file**) baza de date. Am face o operație în plus. Figura 7.11 prezintă opțiunile din **Frequencies** și **Descriptives**. Meniul **Explore** le calculează implicit.

**Figura 7.11.** Meniurile **Frequencies** și **Descriptives**: calcularea skewness și kurtosis





Aici, indicatorii **skewness** și **kurtosis** arată ușoare abateri, dar nimic grav : valori absolute mai mici decât 1 la skewness pentru toate pozițiile sociale și doar o valoare de aproximativ 1.5 la kurtosis pentru partea de sus a clasei mijlocii (tabelul 7.2).

Putem combina valorile de la skewness și kurtosis cu cele de la medie, mediană și medie. Dacă media și mediana au valori apropiate, este mai probabil să nu avem cazuri extreme. Abaterea standard ne va ajuta să înțelegem, de asemenea, cât de omogene sunt grupurile. Să ținem minte însă că am considerat scala de 10 puncte ca fiind metrică. Rezultatele pot fi influențate de acest lucru. De asemenea, să ținem minte că nu toate variabilele au o distribuție naturală normală, deci să nu căutăm normalitate acolo unde este greu de găsit.

**Tabelul 7.2.** Skewness și kurtosis. Calcule efectuate în meniul Explore.  
Tabele obținute prin pivotare

| Descriptives                    |                                   |           |            |
|---------------------------------|-----------------------------------|-----------|------------|
| Statistics= Skewness            |                                   |           |            |
|                                 | V238 Social class<br>(subjective) | Statistic | Std. Error |
| V23 Satisfaction with your life | 1 Upper class                     | -.532     | .553       |
|                                 | 2 Upper middle class              | -.927     | .141       |
|                                 | 3 Lower middle class              | -.641     | .117       |
|                                 | 4 Working class                   | -.477     | .106       |
|                                 | 5 Lower class                     | .030      | .192       |

| Descriptives                    |                                   |           |            |
|---------------------------------|-----------------------------------|-----------|------------|
| Statistics= Kurtosis            |                                   |           |            |
|                                 | V238 Social class<br>(subjective) | Statistic | Std. Error |
| V23 Satisfaction with your life | 1 Upper class                     | -1.038    | 1.069      |
|                                 | 2 Upper middle class              | 1.472     | .281       |
|                                 | 3 Lower middle class              | .446      | .233       |
|                                 | 4 Working class                   | -.452     | .211       |
|                                 | 5 Lower class                     | -.945     | .381       |

## 7.2. Relația liniară dintre două variabile

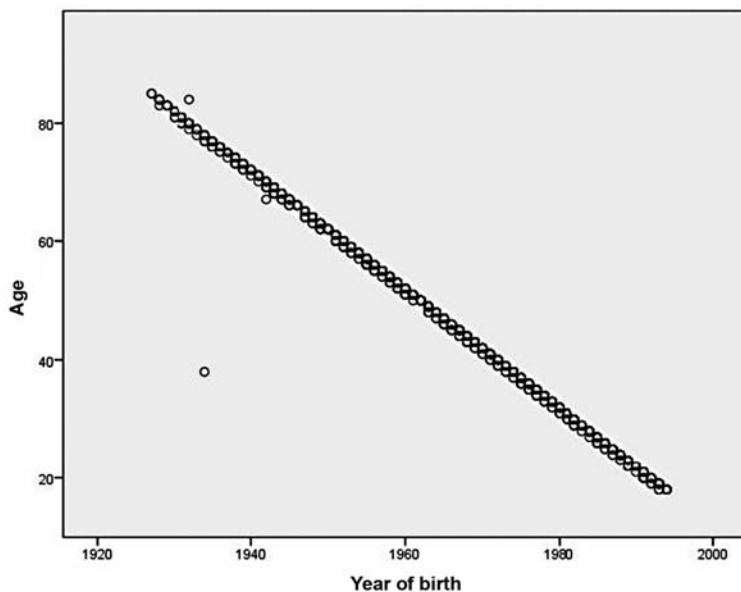
Această asumptie privește relațiile dintre două variabile. Asumptia discutată anterior, chiar dacă am făcut trimiteri și către normalitatea multivariată, a fost dezvoltată într-un cadru univariat.

Analiza de corelație, care calculează coeficientul de corelație Pearson  $r$ , este frecvent utilizată în cercetarea socială. Aceasta ne arată dacă, între două variabile metrice, există o relație: la o anumită valoare a variabilei  $X$ , variabila  $Y$  ia o anumită valoare. Mai general, pentru o mulțime de persoane, dacă valorile variabilei  $X$  cresc sau scad, atunci cresc sau scad și valorile variabilei  $Y$  (relație direct proporțională), sau dacă valorile variabilei  $X$  cresc sau scad, atunci scad sau cresc valorile variabilei  $Y$  (relație invers proporțională). Coeficientul de corelație Pearson  $r$  ia valori în intervalul  $[-1, 1]$ : când  $r = -1$ , relația este perfect negativă, când  $r = 0$ , între  $X$  și  $Y$  nu există o relație, iar când  $r = 1$ , relația este perfect pozitivă. Analiza este simplu de interpretat. Acest lucru o face și atractivă, probabil. Totuși, este foarte ușor să greșim atunci când interpretăm coeficientul de corelație Pearson  $r$  dacă nu am verificat două asumptii esențiale ale acestei analize. Prima asumptie se referă la prezența cazurilor extreme (outlieri). A doua asumptie se referă la existența unei relații liniare între cele două variabile, adică o relație care poate fi aproximată printr-o dreaptă.

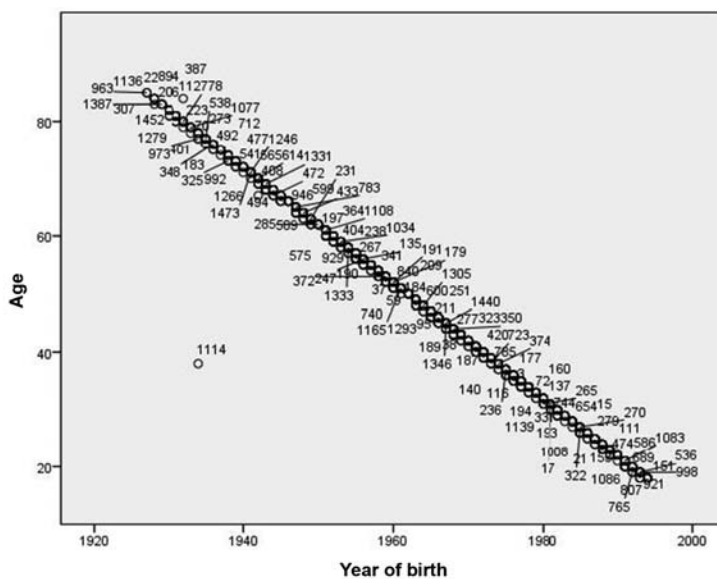
Să luăm relația dintre anul nașterii și vârstă. În baza de date WVS 2012, avem ambele variabile. Vârsta este egală cu anul în care a fost aplicat chestionarul în România, 2012, minus anul nașterii respondentului. Așadar ar trebui să avem o relație perfect liniară negativă: când anul nașterii crește ca valoare, adică este mai apropiat de zilele noastre, vârsta va scădea. Ambele variabile sunt metrice, valorile pe care le pot lua cele două putând fi folosite în calcule aritmetice. Evident, acesta este un exemplu didactic care ne permite să vizualizăm o relație. În practică, nu ar aduce o contribuție prea mare științei investigarea relației dintre acestea. Putem vedea dacă relația este liniară folosind graficul **scatterplot** sau,

în limba română, „nor de puncte”. În figura 7.12 este prezentat acest grafic realizat pentru cele două variabile.

**Figura 7.12.** Scatterplot (nor de puncte) care arată o relație perfect liniară  
(a)



(b)

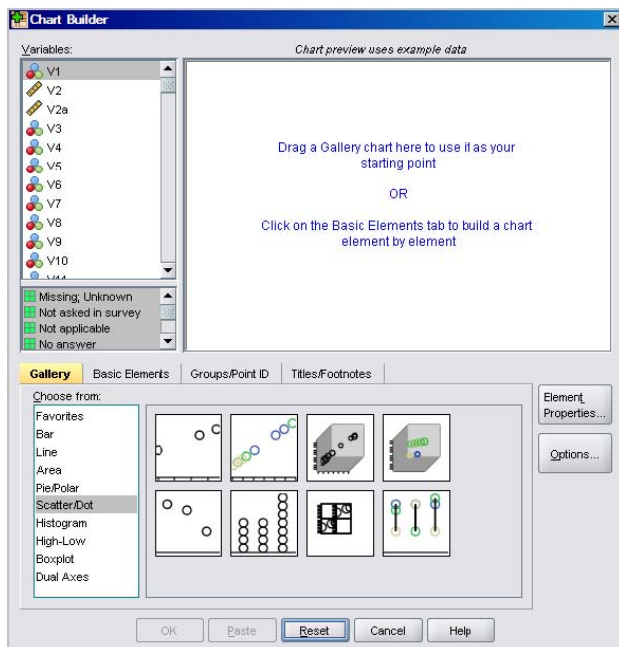


Acest grafic a fost obținut mergând în meniul **Graphs > Chart Builder** (figura 7.13). Puteți lucra și cu **Graphs > Legacy Dialogs** sau **Graphs > Graphboard Template Chooser**. Rezultatele vor fi, în principiu, aceleași.

Așa cum ne-am obișnuit, în stânga, în secțiunea **Variables** sunt toate variabilele din baza de date. De aici vom selecta, pe rând, cele două variabile: anul nașterii (V241) și vârsta în ani împliniți (V242). Imediat sub această fereastră, SPSS afișează valorile variabilei selectate. În partea de jos a ferestrei există patru taburi: **Gallery**, **Basic Elements**, **Groups/Point ID** și **Titles/Footnotes**. În **Gallery** sunt graficele dintre care îl vom alege pe cel care ne interesează. Aici ne interesează scatterplotul, de aceea dăm click pe **Choose from: Scatter/Dot** (figura 7.13a). În partea dreaptă s-au activat opt tipuri de grafice (de la stânga la dreapta): **simple scatter**, **grouped scatter**, **simple 3-d scatter**, **grouped 3-d scatter**, **summary data plot**, **simple dot plot**, **scatterplot matrix** și **drop-line**. Noi vom utiliza graficul **simple scatter**. Mergem cu cursorul pe el și dăm dublu click. Se va activa, în partea centrală a imaginii, structura graficului în care trebuie să introducem informația necesară (figura 7.13b). Pe axa X vom pune variabila pe care o considerăm explicativă. Aici nu are prea mult sens această delimitare între variabilă explicativă (independentă) și variabilă explicată (dependentă). Dar, de dragul prezentării, vom pune anul nașterii (V241) pe axa X, pornind de la ideea că vârsta este derivată din ea. Selectăm V241 și, prin drag-and-drop, o aducem pe axa X. Procedăm similar cu vârsta, V242, dar pe aceasta o ducem pe axa Y.

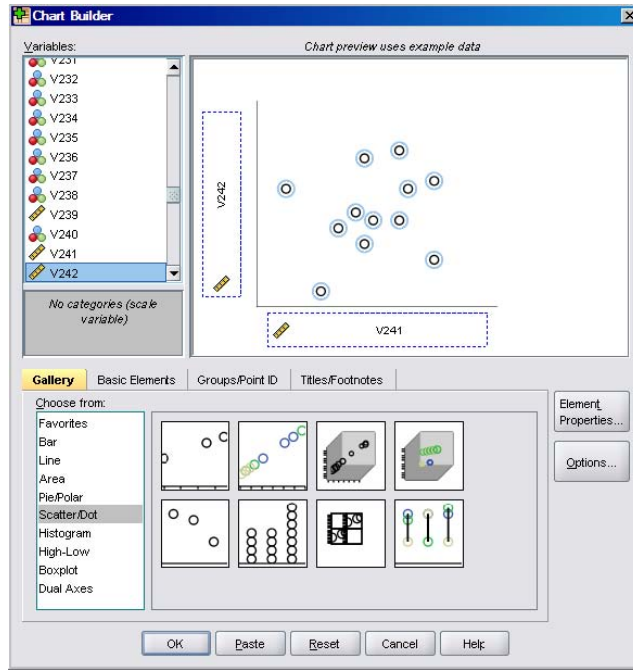
**Figura 7.13.** Meniul Graphs > Chart Builder

(a)

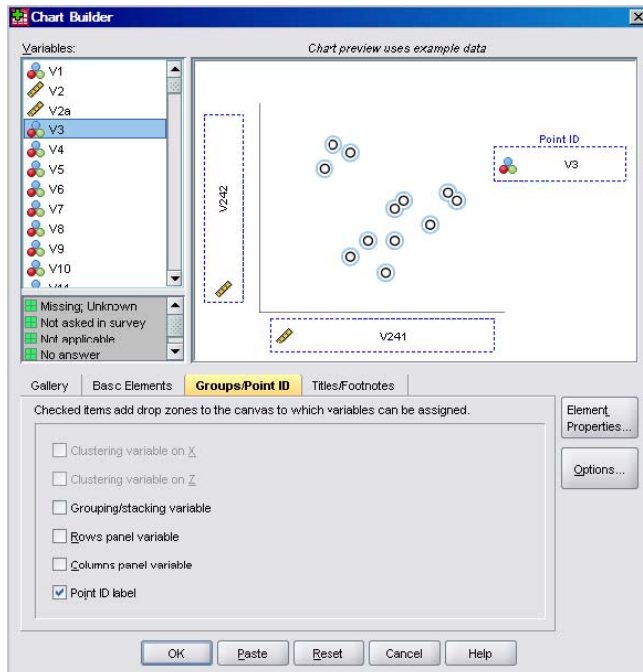




(b)



(c)



Dacă apăsăm **OK**, obținem graficul în care observăm că există un punct ce se abate de la relația așteptată. Pentru a-l putea identifica ușor, ar fi util să știm care este id-ul său unic în baza de date. Nu uitați că, într-o bază de date, toate unitățile introduse trebuie să aibă un id unic. Pentru etichetare trebuie să revenim în meniul **Graphs > Chart Builder** și să utilizăm tabul **Groups/Point ID** (figura 7.13c). Selectăm **Point ID label**. Observăm cum, în fereastra **Chart preview**, a apărut o nouă căsuță intitulată **Point label variable ?**. Aici trebuie să introducem variabila care conține id-urile unice ale respondenților. În baza de date WVS 2012, aceasta este V3. O selectăm și, prin drag-and-drop, o aducem în căsuța activată. Dacă apăsăm **OK**, va rezulta graficul din figura 7.12b. Așadar, persoana care se abate de la relația așteptată are id-ul unic în baza de date 1114. Pentru a vedea ce valori are această persoană la anul nașterii (V241) și vârstă (V242), avem mai multe posibilități. Am putea merge în **Data View**, meniul **Window > Split**. Apoi, am putea căuta valoarea 1114 la V3 fie folosind bara verticală de navigație (scroll), fie folosind procedeul **Find**. Observăm o inadvertență: anul nașterii pentru persoana cu id = 1114 este 1934. În aceste condiții, ne-am aștepta ca vârsta să fie egală cu 78 de ani.

**Figura 7.14.** Window > Split: consultarea vizuală în Data View a unor inadvertențe în date

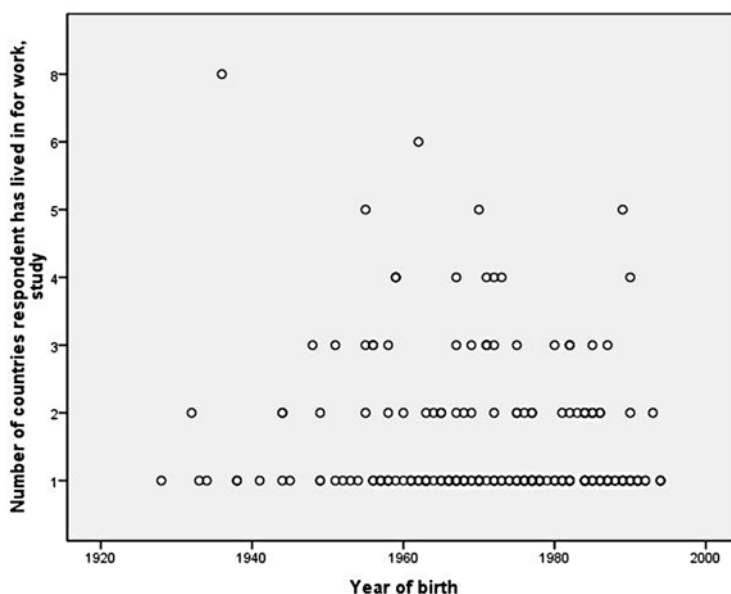
| V241 | V242 |  | V3   |
|------|------|--|------|
| 1934 | 38   |  | 1114 |
| 1955 | 57   |  | 1115 |
| 1993 | 19   |  | 1120 |
| 1977 | 35   |  | 1126 |

Dar aici este egală cu 38 de ani. Dacă nu ar exista astfel de erori, corelația Pearson dintre anul nașterii și vârstă ar fi egală cu  $-1$ . Acum însă este redusă la  $-0.99$ . Aici nu simțim foarte clar efectul cazurilor extreme pentru că, practic, realizăm o corelație a unei variabile cu ea însăși. Însă, atunci când corelăm două variabile care măsoară lucruri diferite, efectul cazului extrem ar putea fi dramatic. O a doua metodă prin care putem vedea ce valori ia cazul cu id-ul 1114 la V241 și V242 presupune următoarele: activăm un filtru care respectă condiția  $V3 = 1114$  și apoi realizăm câte un tabel de frecvență pentru fiecare dintre cele două variabile. Prima metodă este preferată de începătorii care vor să vadă datele. A doua metodă ar trebui însă să fie cea pentru care optăm deoarece ne permite să salvăm sintaxele celor două acțiuni păstrând astfel și jurnalul activității de analiză.

Vizualizarea relațiilor liniare prin utilizarea scatterploturilor este posibilă doar când ambele variabile sunt cu adevărat metrice și au valori multe. Observăm în figura 7.15 un scatterplot cu anul nașterii pe axa X și numărul de țări în care

respondentul a muncit sau studiat. Deși ambele sunt metrice, numărul de țări are o distribuție în formă de J, adică majoritatea persoanelor indică un număr mic de țări. Este greu să apreciem, din acest grafic, ce fel de relație este între cele două variabile.

**Figura 7.15.** Limitele scatterplotului : când o variabilă are puține valori



Jaccard și Jacoby (2010) oferă o explicație frumoasă și ușor de înțeles a funcției liniare. Totuși, acest subiect capătă și mai mult sens dacă se trece într-un cadru multivariat. Regresia liniară multiplă are mai multe asumptii, iar înțelegerea acestora și a metodelor lor de testare va face mult mai clar modul în care putem depista relații nonliniare între variabilele care ne interesează (Berry, 1993).

### 7.3. Soluții la încălcarea asumptiei de normalitate a distribuției

Deseori, asumptia normalității este încălcată datorită prezenței cazurilor extreme. Uneori, de exemplu, când sunt foarte puține, putem șterge cazurile extreme, rezolvând astfel și problema normalității. Alteori însă lucrurile nu sunt atât de simple.

O soluție des utilizată de cercetătorii experimentați, atunci când asumptia de normalitate este încălcată, constă în transformarea matematică a uneia dintre

variabile sau chiar a mai multora. Hair *et al.* (2010), Tabachnick și Fidell (2007) sau Field (2007) sunt doar câteva dintre lucrările în care ne sunt prezentate alternativele pe care le avem la îndemână. Trebuie însă să reținem că aceste transformări nu sunt soluții minune. De exemplu, interpretarea modelelor în care sunt folosite variabile transformate este mai dificilă decât atunci când utilizăm unitățile de măsură originale.

Pentru că depășește scopul acestei lucrări, vom reda doar câteva dintre transformările uzuale, așa cum sunt recomandate de autorii citați :

- distribuție alungită la dreapta, distanța dintre valorile minime și maxime „normale” este mică : logaritmăm variabila ;
- distribuție alungită la dreapta, distanța dintre valorile minime și maxime „normale” este ceva mai mare : radical din variabilă ;
- distribuție alungită la dreapta, cu formă care aproximează litera J întoarsă (censored) : calculăm raportul dintre 1 și variabilă ( $1/\text{variabilă}$ ).

O lucrare foarte utilă pentru cei care vor să pătrundă tainele acestor probleme matematice îi aparține lui John Fox (2009).

## 7.4. Exerciții

Pentru aceste exerciții utilizăm baza de date și/sau chestionarul World Values Survey 2012 rezultată(e) în urma aplicării chestionarului în România. Baza de date poate fi descărcată de pe pagina de internet a *Grupului Românesc pentru Studiul Valorilor Sociale* (<http://www.romanianvalues.ro>).

1. Verificați asumpțiile pentru toate exercițiile de la capitolul 6.
2. Propuneți soluții de îmbunătățire a situației acolo unde este cazul.

## 8. Corelația și regresia liniară multiplă

Care este relația dintre veniturile unei persoane și numărul anilor de educație formală absolviți? Are sens investiția de timp și resurse în educație? Cresc veniturile odată cu numărul anilor de educație formală absolviți? Notele primite la testul-grilă cresc odată cu numărul de cursuri și seminarii la care studenții participă? Sau, mai degrabă, notele tind să fie mai mari atunci când studenții petrec mai multe ore studiind individual? Satisfacția cu viața este mai ridicată atunci când persoanele consideră că au control asupra propriei vieți? Acesta este un tip de întrebări pe care ni le punem frecvent în procesul de cercetare. De fapt, abstractizând, ne întrebăm dacă între două variabile există o corelație. Termenul corelație este nou, dar ideea nu, aceasta devenind familiară deja de la asocierea testată prin tabele de contingență și chi square (hi pătrat).

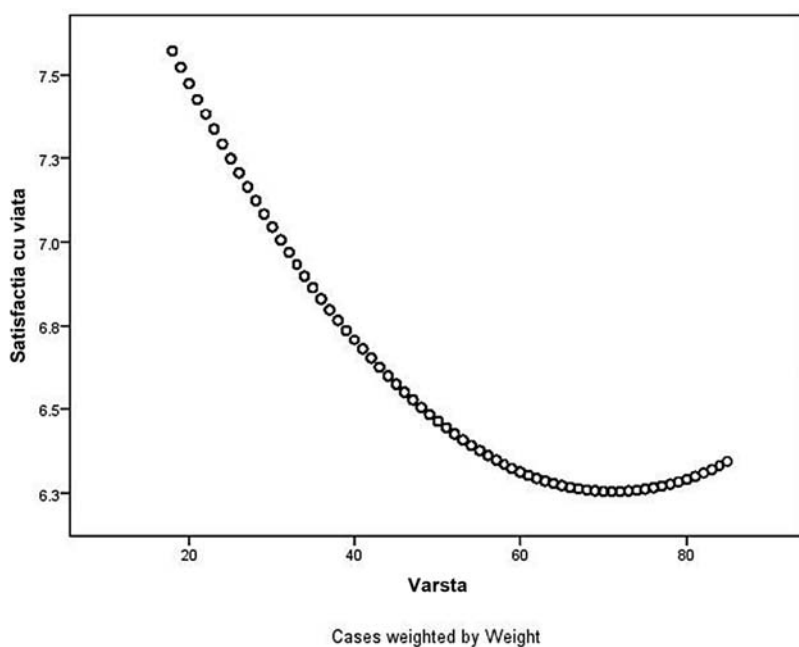
Acesta este un pas intermediar către analiza multivariată. De fapt, viața socială este complexă și nu poate fi explicată adecvat folosind analize bivariate. Trebuie să învățăm cum putem utiliza, simultan, mai mulți predictorii pentru aceeași variabilă dependentă. Veniturile unei persoane depind atât de educația formală acumulată, cât și de capitalul cultural moștenit de la părinți, generația din care face parte, sexul acesteia, vârsta, sistemul de valori la care aderă, tipul localității în care locuiește etc. Notele primite la testul-grilă depind de prezența fizică la ore deoarece studentul poate nota idei și exemple oferite spontan de profesor, poate pune întrebări prin care să își clarifice concepte și moduri de lucru, dar și pentru că poate discuta în pauze cu colegii despre ceea ce au auzit în timpul întâlnirii. De asemenea, contează și studiul individual. Dar toate aceste relații se pot modifica dacă luăm în considerare, de exemplu, gradul de extraversiune al studenților. O persoană care consideră că are control asupra propriei vieți ar putea fi mai mulțumită cu viața sa. Dar mulțumirea depinde și de starea de sănătate, de starea civilă, situația financiară etc. Așadar, o explicație cere considerarea simultană a mai multor factori care determină variația fenomenului care ne interesează. Una dintre cele mai utilizate analize multivariate este regresia liniară multiplă.

În prima parte a capitolului vom discuta despre corelația liniară necesară pentru înțelegerea regresiei liniare multiple. Apoi, vom prezenta, într-o manieră nontehnică, principiile regresiei liniare multiple și aplicabilitatea acesteia.

## 8.1. Corelația liniară

Numele analizei bivariante discutate aici implică asumția că, între cele două variabile pentru care calculăm coeficientul de corelație, există o relație liniară, adică o relație care poate fi reprezentată printr-o dreaptă. Relațiile sunt de tipul: (a) X1 crește, X2 crește; (b) X1 crește, X2 scade; (c) X1 scade, X2 crește; (d) X1 scade, X2 scade. Nu există puncte de inflexiune. Atunci când există puncte de inflexiune, relația nu mai este liniară. De exemplu, relația dintre vârstă și satisfacția cu viața nu este liniară. Pentru a înțelege mai bine, să privim figura 8.1.

**Figura 8.1.** Relația nonliniară dintre vârstă și satisfacția cu viața



Satisfacția cu viața are valoare maximă începând cu 18 ani (vârsta minimă în eșantionul WVS 2012 pentru România). Aceasta descrește constant, dar, la o anumită vârstă, pare să revină pe un trend ascendent. Acel punct de inflexiune arată o relație nonliniară între cele două variabile.

O altă asumție a corelației liniare este că ambele variabile sunt cantitative, adică, în termenii nivelurilor de măsurare, de interval sau raport. În științele sociale, frecvent, scalele simple tip Likert (o întrebare cu minim 4 variante de răspuns de tipul acord/dezacord), dar și scorurile derivate din cele compuse sunt considerate de interval, deci cantitative. De exemplu, satisfacția cu viața măsurată

pe o scală de la 1 la 10 este utilizată deseori în analize în acest mod. Nu are sens calculul coeficientului de corelație pentru variabilele categoriale, adică nominale și ordinale de tipul categorii ordonate.

Corelația nu înseamnă cauzalitate. Analiza de corelație ne arată doar că două variabile variază împreună, felul relației (direct sau invers proporțională) și cât de puternică este aceasta. Însă nu putem spune cu certitudine că  $X_1$  o determină pe  $X_2$  sau invers. Pentru o interpretare în termeni cauzali, cercetătorul trebuie să respecte o serie de principii chiar în designul cercetării, cum ar fi opțiunea pentru experiment sau anchetă prin chestionar. În practică, cercetătorul nu gândește în termenii  $X_1$  și  $X_2$ , ci în termenii  $X$  și  $Y$ , adică o variabilă independentă și una dependentă. Adecvarea și consistența interpretării ține de corectitudinea logicii cercetătorului. Un domeniu în care erorile de interpretare ale analizei de corelație sunt foarte posibile este cel al fericirii și satisfacției cu viața. Care este, spre exemplu, relația dintre starea de sănătate percepută și satisfacția cu viața? O persoană despre care se consideră că este mai sănătoasă va fi mai satisfăcută cu viața sau o persoană mai satisfăcută cu viața se va considera mai sănătoasă? Răspunsul nu este unul simplu, ambele variante având un anumit grad de plauzibilitate. În multe studii însă, satisfacția cu viața este considerată variabila dependentă, iar starea subiectivă de sănătate este considerată variabila independentă. Starea subiectivă de sănătate depinde, în mare măsură, de starea obiectivă de sănătate (prezența unei boli temporare, a unei boli cronice, a unui handicap etc.), deci, dacă se intervine asupra stării obiective de sănătate, se va ajusta și starea subiectivă de sănătate și, într-un final, satisfacția cu viața, privită ca un rezultat al vieții de calitate. Iar o viață de calitate cu o stare de sănătate precară este destul de greu de imaginat.

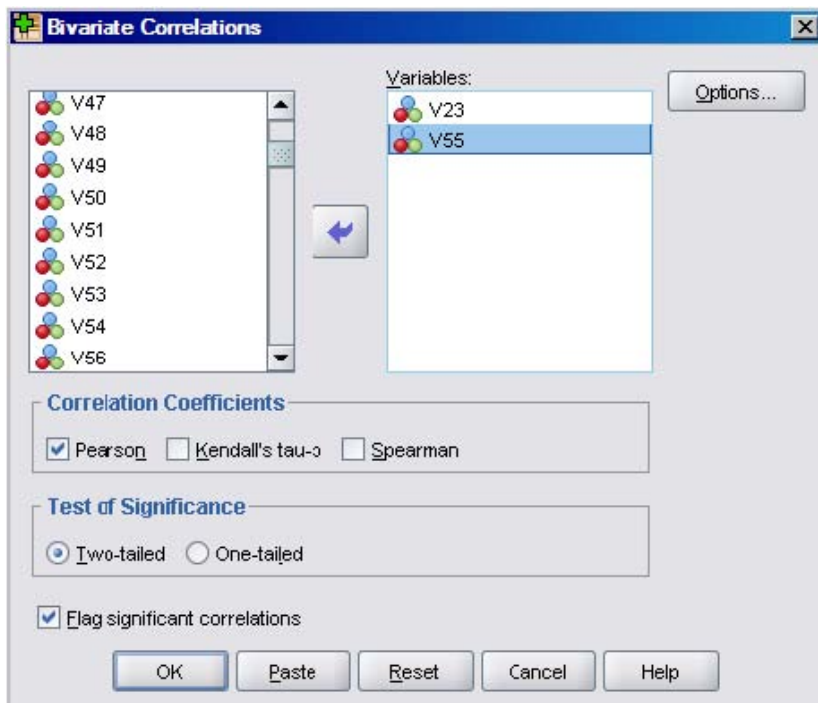
Există mai mulți indicatori de corelație. Alegerea unuia depinde de caracteristicile variabilelor pe care dorim să le corelăm: variabilele corelate pot să ia multe valori, iar acestea sunt numere; variabilele corelate conțin ranguri naturale; variabilele corelate sunt dihotomice naturale sau sunt dihotomice obținute prin recodificarea unor variabile continue etc. Aici discutăm despre coeficientul de corelație Pearson  $r$ , care presupune că ambele variabile sunt cantitative continue (valorile variabilelor sunt numerice și destul de multe). Alături de acesta, mai des întâlniți în practică sunt coeficienții de corelație Spearman, Kendall și Gamma. Aceștia sunt specifici variabilelor ordinale, dar sunt utilizați și atunci când anumite asumptii, cum ar fi cea a distribuției normale bivariate, sunt încălcate. Sunt corelații nonparametrice, spre deosebire de Pearson, care este parametrică. O descriere foarte bună a celor mai utilizați coeficienți de corelație a fost realizată de Chen și Popovich (2002).

În SPSS analiza de corelație are un meniu dedicat: **Analyze > Correlate**. SPSS poate calcula două tipuri de corelații din acest meniu: corelație bivariată și corelație parțială. În același meniu mai există opțiunea calculării unor distanțe între cazuri sau variabile. Această analiză este însă în afara intereselor noastre și

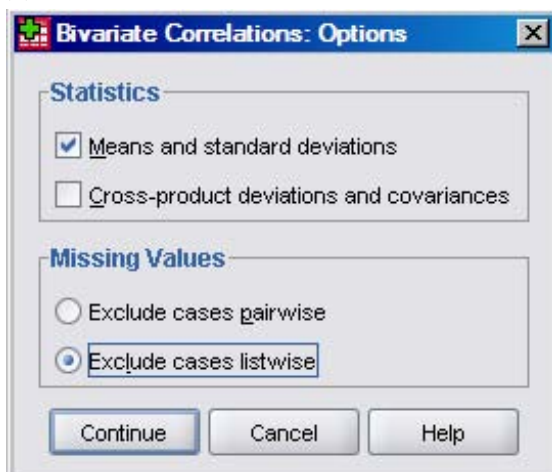
poate fi înțeleasă mai bine în contextul analizei cluster. Meniul pentru corelația bivariată, **Analyze > Correlate > Bivariate**, ne permite calcularea a trei coeficienți de corelație : **Pearson**, **Kendall's tau-b** și **Spearman** (figura 18.2).

**Figura 8.2.** Meniul corelației bivariate (Correlate > Bivariate)

(a)



(b)





În secțiunea din stânga sus sunt toate variabilele din baza de date. De aici selectăm variabilele pe care dorim să le corelăm și, folosind săgeata dintre cele două secțiuni, le trecem în secțiunea din dreapta sus. Trebuie să introducem minim două variabile. Nu este obligatoriu să introduceți doar două. Dacă introducem opt variabile pentru care dorim să calculăm coeficientul de corelație Pearson, atunci va rezulta un tabel cu opt rânduri și opt coloane, fiecare dintre cele opt variabile fiind corelată, pe rând, cu ea însăși și cu toate celelalte șapte. Așadar, indiferent de numărul de variabile pe care le introducem în analiza corelației, rezultatul va fi tot bivariat. În figura 8.2 am introdus doar două variabile, v23 și v55. Prima reflectă răspunsurile la întrebarea din WVS 2012 „Dacă luați în considerare toate aspectele vieții dvs. din ultimul timp, în ce măsură sunteți mulțumit de ea ? (utilizați scala de mai jos, în care 1 înseamnă «total nemulțumit» și 10 «total mulțumit»)”. A doua reflectă răspunsurile la întrebarea din aceeași cercetare „Unii oameni cred că au libertate totală de alegere și de control asupra vieții lor, iar alți oameni cred că, indiferent ce fac, nu pot influența ce li se întâmplă în viață. Vă rugăm să folosiți scala următoare pentru a indica câtă libertate de alegere credeți că aveți dvs., dând o notă de la 1 la 10, unde 1 înseamnă că «Nu am deloc», iar 10 că «Am libertate deplină»”.

Următorul pas constă în alegerea coeficientului de corelație pe care dorim să-l calculăm. Interpretarea coeficientului de corelație este relativ simplă și directă. Dacă valoarea  $p$  calculată este mai mică sau egală cu pragul teoretic de 0.05, atunci consultăm semnul coeficientului, care ne spune direcția relației, urmând ca mai apoi să interpretăm puterea relației dată de valoarea absolută a coeficientului de corelație.

Toți cei trei coeficienți, Pearson, Kendall tau-b și Spearman, variază între  $[-1, 1]$ .

Interpretarea semnului se face în funcție de semnificația valorilor pe care le iau cele două variabile analizate. Dacă variabilele sunt numere, atunci interpretarea este simplă. Când numărul de ore petrecute studiind individual la statistică crește, ne așteptăm ca notele luate la teste să crească. Semnul va fi plus. Când numărul anilor de educație formală crește, ne așteptăm ca veniturile persoanei să crească. Semnul va fi plus. Când numărul sortimentelor vândute într-un magazin este mare, ne așteptăm ca numărul clienților aceluși magazin să fie mare. Semnul va fi plus. Dar dacă calculăm coeficientul de corelație Pearson pentru două variabile măsurate fiecare pe o scală de 10 puncte, atunci trebuie să citim cu atenție etichetele atribuite codurilor. Dacă 1 înseamnă satisfacție scăzută cu viața și 10 satisfacție ridicată, iar 1 înseamnă absența controlului asupra propriei vieți și 10 înseamnă control total, atunci semnul va fi plus. Dacă una dintre cele două variabile ar fi codificată în alt sens, de exemplu 1 ar însemna satisfacție ridicată cu viața și 10 satisfacție scăzută, iar la control scala s-ar păstra, atunci semnul ar fi minus.

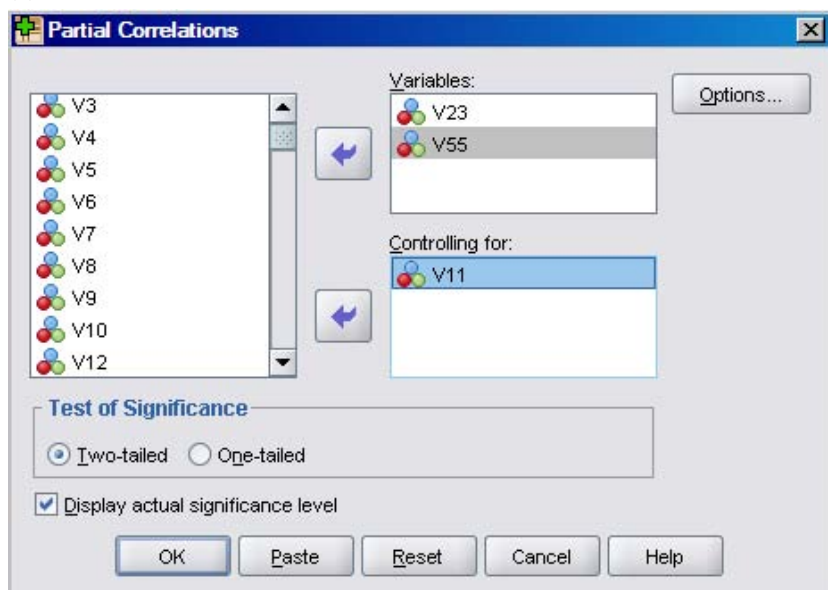
Valorile extreme indică relații perfecte de interdependență între cele două variabile. Valoarea 0 poate însemna absența unei relații de interdependență. Dar,

dacă ne reamintim de asumptia liniarității, dar și de cea a distribuției bivariate normale, un coeficient Pearson egal cu zero poate sugera și o relație nonliniară, atunci când aceste asumptii nu sunt respectate. Adică între cele două variabile există o relație care nu este liniară (figura 8.1), deci nu poate fi reprezentată numeric corect prin coeficientul Pearson. Nu există o regulă larg acceptată despre relația dintre valoarea coeficientului de corelație și tăria corelației. De regulă, valorile absolute mai mici de 0.3 sunt considerate corelații slabe spre moderate, între 0.3 și mai mici de 0.6 sunt considerate corelații moderate spre puternice, iar mai mari sau egale cu 0.6 sunt considerate corelații puternice. Totuși, interpretarea depinde de domeniul studiat, iar aceste valori au un caracter orientativ. Trebuie ținut cont și de forma distribuțiilor celor două variabile: când sunt diferite de cea normală și/sau diferite între ele, atunci valorile maxime,  $-1$  sau  $1$ , sunt mai greu de atins (Carroll, 1961). De asemenea, Chen și Popovich (2002) atrag atenția că în eșantioanele mici, de câteva zeci de cazuri, este foarte probabil să avem coeficienți de corelație cu valori mari, chiar dacă în populație valorile sunt mici sau corelația este inexistentă. În aceeași situație, trebuie să fim atenți și la cazurile extreme care pot afecta mărimea sau chiar direcția coeficientului de corelație Pearson. Lucrarea celor doi autori detaliază toate problemele care pot influența rezultatul analizei de corelație, atunci când folosim coeficientul de corelație Pearson.

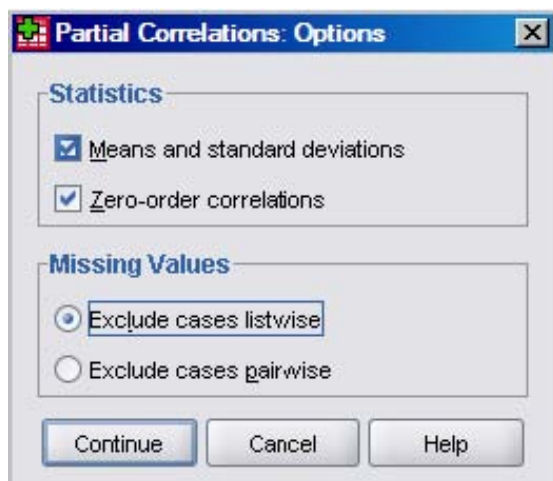
O altă problemă pe care trebuie să o avem în vedere când interpretăm un coeficient de corelație Pearson se referă la relația iluzorie dintre cele două variabile corelate. Relația iluzorie (*spurious*) atestă existența a cel puțin unei alte variabile,  $X_3$ , care explică relația dintre  $X_1$  și  $X_2$ . Kline (2011) arată cum relația dintre numărul de cuvinte pe care un copil îl are în vocabular ( $X_1$ ) și mărimea pe care o poartă la pantof ( $X_2$ ) este iluzorie, ambele fiind, de fapt, determinate de vârsta copilului ( $X_3$ ) folosită aici ca proxy pentru maturizarea educațională și fizică. Astfel de situații trebuie identificate teoretic de cercetător în acord cu literatura consultată sau, atunci când aceasta lipsește, cu intuiția proprie. Practic, poate fi testată folosind corelația parțială. Însă, deși pare atrăgătoare această metodă de testare a veridicității relației dintre două variabile, în realitate este insuficientă. Viața socială este mult mai complexă, fiind foarte probabil să existe mai mult de o variabilă care explică relația dintre cele două variabile corelate. Aceste situații pot fi testate mai adecvat în contextul modelelor de ecuații structurale (*structural equation modeling*), un subiect care depășește tematica acestui volum. În SPSS, corelația parțială poate fi găsită în meniul **Analyze > Correlate > Partial** (figura 8.3).

Figura 8.3. Meniul corelației parțiale (Correlate &gt; Partial)

(a)



(b)



În secțiunea din stânga a ferestrei, care se deschide prin activarea meniului, se găsesc variabilele din baza de date. În secțiunea denumită **Variables**, introducem cele două variabile pe care dorim să le controlăm. Aici am introdus V23, satisfacția cu viața, și V55, controlul perceput asupra propriei vieți. În secțiunea denumită **Controlling for** introducem variabila despre care presupunem că influențează relația dintre V23 și V55. Aici am introdus V11, care cuprinde răspunsurile la

întrebarea din WVS 2012 : „Cum ați descrie starea dvs. de sănătate în prezent ? 1 Foarte bună, 2 Bună, 3 Nu prea bună, 4 Proastă”. Presupunem că persoanele care consideră că au o stare de sănătate pozitivă apreciază că au un control ridicat asupra propriei vieți. Dacă se întâmplă acest lucru, atunci ne așteptăm ca relația dintre controlul perceput asupra propriei vieți și satisfacția cu viața să se diminueze, adică valoarea coeficientului de corelație Pearson să scadă. Rughiniș (2007) dă mai multe exemple de corelații iluzorii, explicând care este rolul controlului variabilelor în analiză. Astfel, putem explica relația dintre numărul de pompieri la locul unui incendiu și valoarea pagubelor produse dacă luăm în considerare mărimea incendiului. La fel, putem înțelege relația dintre numărul bisericilor dintr-un oraș și numărul crimelor violente, ambele fiind determinate de mărimea localității.

În meniul corelației bivariante (figura 8.2), dar și în cel al corelației parțiale (figura 8.3), mai avem câteva opțiuni pe care le putem bifa sau debifa.

Una dintre ele se referă la tipul testului de semnificație: **two-tailed** sau **one-tailed**. Opțiunea implicită în SPSS este two-tailed. Cercetătorul poate păstra sau modifica această opțiune în funcție de ipoteza pe care o testează. Pentru o înțelegere a logicii testelor de semnificație, poate fi consultat orice manual de statistică sau lucrările dedicate acestui subiect de către Henkel (1976) și Mohr (1990). În practică, așa cum subliniază Mohr (1990), majoritatea ipotezelor sunt direcționale, adică presupunem că relația are o anumită direcție. Din acest motiv, ar trebui să selectăm one-tailed. Însă același autor consideră că mulți cercetători adoptă o perspectivă mai conservatoare și aleg opțiunea two-tailed, ca și când natura relației nu ar putea fi precisă, folosind termenii lui Field (2009). De exemplu, dacă credem că persoanele care consideră că au control asupra propriei vieți vor fi mai satisfăcute cu viața lor, atunci aș putea alege opțiunea one-tailed. Dar dacă nu știm la ce să ne așteptăm, alegem opțiunea two-tailed. E preferabil să alegem varianta conservatoare și să păstrăm opțiunea implicită din SPSS.

O altă opțiune este **Flag significant correlations** în meniul corelației bivariante (figura 8.2) și **Display actual significance level** în meniul corelației parțiale (figura 8.3). Ambele opțiuni au efect doar asupra modului de prezentare a tabelelor în Output. Prefer opțiunea implicită din SPSS. În tabelul 8.1 sunt prezentate rezultatele cu și fără aceste opțiuni bifate. În cazul corelației bivariante, când păstrăm bifată opțiunea implicită, în dreptul coeficientului de corelație sunt notate una sau mai multe steluțe (\*, \*\*), în funcție de valoarea pe care o ia nivelul de semnificație. Când lucrăm, acest ajutor vizual poate fi foarte util, de aceea recomand utilizarea sa. În cazul corelației parțiale, lucrurile stau invers: prin debifare sunt afișate steluțele în defavoarea valorii nivelului de semnificație. Tabelul afișat este mai puțin complex, dar dacă în lucrarea pe care o pregătim trebuie să raportăm chiar nivelul de semnificație calculat, atunci am avea nevoie să păstrăm bifată opțiunea implicită.

**Tabelul 8.1.** Output cu sau fără opțiunile Flag... sau Display... în meniurile corelației bivariate, respectiv corelației parțiale

| <b>Correlations</b>  |                     |                                       |  |
|--|---------------------|---------------------------------------|--|
|  |                     | V23<br>Satisfaction<br>with your life | V55 How much<br>freedom of<br>choice and<br>control over<br>own life |
| V23 Satisfaction with your<br>life                             | Pearson Correlation | 1                                     | .333**   |
|  | Sig. (2-tailed)     |                                       | .000   |
|  | N                   | 1491                                  | 1474   |
| V55 How much freedom of<br>choice and control over<br>own life | Pearson Correlation | .333**                                | 1  |
|  | Sig. (2-tailed)     | .000                                  |  |
|  | N                   | 1474                                  | 1484   |
| **. Correlation is significant at the 0.01 level (2-tailed).   |                     |                                       |  |

| <b>Correlations</b>   |                     |                                       |  |
|---|---------------------|---------------------------------------|--|
|   |                     | V23<br>Satisfaction<br>with your life | V55 How much<br>freedom of<br>choice and<br>control over<br>own life |
| V23 Satisfaction with your<br>life                          | Pearson Correlation | 1                                     | .333   |
|   | Sig. (2-tailed)     |                                       | .000   |
|   | N                   | 1491                                  | 1474   |
| V55 How much freedom of<br>choice and control over own life | Pearson Correlation | .333                                  | 1  |
|   | Sig. (2-tailed)     | .000                                  |  |
|   | N                   | 1474                                  | 1484   |

| <b>Correlations</b>                 |  |                            |                                       |  |
|-------------------------------------|--|----------------------------|---------------------------------------|--|
| Control Variables                   |  |                            | V23<br>Satisfaction<br>with your life | V55 How much<br>freedom of<br>choice and<br>control over<br>own life |
| V11 State of<br>health (subjective) | V23<br>Satisfaction<br>with your life                                | Correlation                | 1.000                                 | .301   |
|                                     |  | Significance<br>(2-tailed) | .                                     | .000   |
|                                     |  | df                         | 0                                     | 1471   |
|                                     | V55 How<br>much freedom<br>of choice and<br>control over<br>own life | Correlation                | .301                                  | 1.000  |
|                                     |  | Significance<br>(2-tailed) | .000                                  | .  |
|                                     |  | df                         | 1471                                  | 0  |

| Correlations                                 |   |             |  |   |
|--|---|-------------|--|---|
| Control Variables                            |   |             | V23<br>Satisfaction<br>with your<br>life | V55 How<br>much<br>freedom of<br>choice and<br>control over<br>own life |
| V11 State of health<br>(subjective)          | V23 Satisfaction with<br>your life                                | Correlation | 1.000                                    | .301**  |
|  | V55 How much<br>freedom of choice<br>and control over own<br>life | Correlation | .301**                                   | 1.000   |
| **. Correlation is significant at 0.01 level |   |             |  |   |

În ambele meniuri, există butonul **Options** care activează opțiunile prezentate în figurile 8.2b și 8.3b. Cu excepția, **Cross-product deviations and covariances** și **Zero-order correlations**, celelalte opțiuni sunt similare. În practică, de regulă, în cazul corelației bivariante, bifăm **Means and standard deviations**, iar în cazul corelației parțiale bifăm această opțiune și **Zero-order correlations**. Prima opțiune ne afișează media și abaterea standard pentru fiecare dintre variabilele incluse în analiză (tabelul 8.2).

**Tabelul 8.2.** Opțiunea Means and standard deviations din meniurile corelației bivariante, respectiv corelației parțiale

| Descriptive Statistics   |      |                |      |
|--|------|----------------|------|
|  | Mean | Std. Deviation | N    |
| V23 Satisfaction with your life                                | 6.70 | 2.385          | 1474 |
| V55 How much freedom of<br>choice and control over own<br>life | 7.88 | 2.279          | 1474 |

Faptul că ni se oferă posibilitatea de a calcula media și abaterea standard pentru variabilele corelate face evident, încă o dată, că această analiză solicită variabile metrice continue. În exemplul prezentat aici, am utilizat două scale simple tip Likert cu 10 variante de răspuns și o scală simplă tip Likert cu 4 variante de răspuns. Mediile și, implicit, abaterile standard, calculate pentru astfel de variabile au un caracter mai degrabă artificial fiind, uneori, chiar dificil de interpretat (de exemplu, când scala are o variantă de mijloc evidențiată printr-o etichetă de tipul „nici acord, nici dezacord”). Chiar dacă în practică astfel de analize sunt acceptate convențional, trebuie să fim conștienți de posibilele erori pe care le putem introduce în interpretările substanțiale ale unor astfel de rezultate.

**Tabelul 8.3.** Opțiunea Zero-order correlations în meniul corelației parțiale

| <b>Correlations</b>                                 |   |             |  |   |   |
|---|---|-------------|--|---|---|
| Control Variables                                   |   |             | V23<br>Satisfac-<br>tion with<br>your life | V55 How<br>much<br>freedom<br>of choice<br>and<br>control<br>over own<br>life | V11 State<br>of health<br>(subjec-<br>tive) |
| -none- <sup>a</sup>                                 | V23 Satisfaction<br>with your life                                | Correlation | 1.000                                      | .333**  | -.365**                                     |
|   | V55 How much<br>freedom of choice<br>and control over<br>own life | Correlation | .333**                                     | 1.000   | -.154**                                     |
|   | V11 State of<br>health (subjective)                               | Correlation | -.365**                                    | -.154**   | 1.000                                       |
| V11 State of<br>health (subjec-<br>tive)            | V23 Satisfaction<br>with your life                                | Correlation | 1.000                                      | .301**  |   |
|   | V55 How much<br>freedom of choice<br>and control over<br>own life | Correlation | .301**                                     | 1.000   |   |
| a. Cells contain zero-order (Pearson) correlations. |   |             |  |   |   |
| **. Correlation is significant at 0.01 level        |   |             |  |   |   |

Zero-order correlations (tabelul 8.3) se referă la corelațiile bivariante dintre toate variabilele pe care le includem în analiza de corelație parțială. Aici avem trei variabile: v23, v55 și v11. Corelația parțială cu o singură variabilă de control se numește first-order correlation.

În fine, ultimul lucru care ne interesează, la acest nivel, este modul de tratare a nonrăspunsurilor în analiza de corelație. Deși nu am menționat până acum, bănuiesc că a fost evident că nu putem calcula coeficientul de corelație Pearson sau oricare altul decât după ce am instruit SPSS să dezactiveze în analize codurile de nonrăspuns (**missing**). Pentru cele trei variabile utilizate pentru exemplificare, tabelele de frecvențe (**Analyze > Descriptive statistics > Frequencies**) arată următoarele :

| Variabila  | Volumul<br>eșantionului | Nonrăspunsuri | Volumul<br>valid |
|--|-------------------------|---------------|------------------|
| V23, satisfacția cu viața                        | 1503                    | 13 / 1 %      | 1490             |
| V55, controlul perceput asupra<br>propriei vieți | 1503                    | 20 / 1 %      | 1483             |
| V11, evaluarea sănătății proprii                 | 1503                    | 1 / 0.1 %     | 1502             |

Nonrăspunsurile (codurile -2 – „nu răspund”, respectiv -1 – „nu știu”) au fost scoase din analiză în **Variable View > Missing > Discrete missing values**. După această operațiune, au fost calculați coeficienții de corelație. Dacă ne uităm în tabelul 8.2, observăm în dreptul celor trei variabile același total, 1.474 persoane. S-a ajuns la același total selectând **Exclude cases listwise** în secțiunea **Missing Values**. Au fost ignorați în analiza de corelație indivizii care nu au oferit un răspuns valid la cel puțin una dintre cele trei variabile analizate. Dacă selectam **Exclude cases pairwise**, atunci am fi avut totaluri diferite la variabile astfel : la corelația bivariată 1.490 la v23 și 1.483 la v55. În practică, pentru a nu introduce erori de interpretare, dată fiind compoziția diferită a grupurilor, alegem să tratăm nonrăspunsurile **listwise**.

În încheiere, să interpretăm și coeficienții de corelație bivariată, respectiv cei de corelație parțială. Folosim informațiile din figura 8.4.

În ambele situații, fie că avem o corelație bivariată, fie una parțială, mai întâi consultăm valoarea p (nivelul de semnificație), care în SPSS este notată **Sig** sau „**Significance**”. Pragul de semnificație este ales a priori analizei. Pragurile acceptate sunt 0.05 și 0.01. Dacă am ales pragul cel mai puțin restrictiv, 0.05, și observăm că p calculat este mai mic decât această valoare, atunci putem respinge ipoteza de nul a lipsei de corelație (am folosit varianta **two-tailed**, non-direcțională). Pentru corelația dintre v23 (satisfacția cu viața) și v55 (controlul perceput asupra propriei vieți) p este egal cu 0.000. De fapt, p nu este 0, ci o valoare cu foarte multe zecimale după virgulă. Fiind mai mică decât 0.05 putem aprecia că există o corelație între satisfacție și control.

Semnul coeficientului de corelație Pearson este „+”, deci am fi tentați să spunem că ambele variabile variază în același sens. Deoarece acestea sunt măsurate prin scale cu 10 puncte, trebuie să vedem cum sunt codificate. Aici sensul este același : codul cel mai mic înseamnă situația negativă (satisfacție scăzută, respectiv lipsa controlului), iar codul cel mai mare înseamnă situația pozitivă (satisfacție ridicată, respectiv prezența unui control ridicat). Așadar semnul pozitiv indică o relație pozitivă. Ne amintim că, teoretic, corelația nu implică cauzalitate. În practică însă, cercetătorul atribuie unei variabile rolul de dependentă, iar celelalte de independentă. Aici am interpreta că, atunci când sentimentul de control asupra propriei vieți crește, crește și satisfacția cu viața.

În fine, trebuie să apreciem cât de puternică este corelația :  $r = 0.33$ . Folosind regulile empirice întâlnite în multe surse științifice, aceasta este o corelație moderată.

Interpretarea este similară pentru coeficientul de corelație parțială. Apare ceva în plus din punct de vedere conceptual, lucru evident în tabelul 8.3. Să ne reamintim că am controlat pentru v11, evaluarea stării de sănătate, pentru că am presupus că explică parțial relația dintre controlul perceput asupra vieții și satisfacția cu viața. Dacă se întâmplă așa, atunci ne așteptăm ca, după ce am controlat pentru v11, corelația dintre v23 și v55 să scadă. Coeficientul de corelație bivariată dintre v23



și  $v_{55}$  este egal cu 0.33, iar după ce am controlat pentru  $v_{11}$  acesta scade la 0.30. Diferența nu este mare. Putem fi entuziaști și să observăm scăderea, dar trebuie să fim și realiști văzând că diferența nu este mare. Probabil mai există și alți factori care modelează relația dintre  $v_{23}$  și  $v_{55}$ . Dar testarea acestei idei presupune un cadru multivariat.

## 8.2. Regresia liniară multiplă

Calitatea vieții unei persoane sau, în ansamblu, a unei populații are două componente: una obiectivă, de stare, și una subiectivă, de evaluare (Zamfir *et al.*, 1984). Aceste două componente presupun efectuarea unor măsurători pe mai multe dimensiuni ale vieții. La nivel individual starea sănătății se poate măsura, printre altele, prin prezența/absența unei boli cronice și/sau a unei incapacități fizice care împiedică persoana, într-o anumită măsură, să își desfășoare activitățile într-o zi obișnuită. La nivel național, starea sănătății se poate măsura, printre altele, folosind rata de morbiditate. Elaborarea unui set comprehensiv de indicatori pentru care poate fi culeasă informație statistică de calitate este dificilă dată fiind complexitatea dimensiunilor vieții umane. O încercare de sistematizare este oferită de Mărginean (2005). Această perspectivă asupra calității vieții ia prea puțin în considerare persoana ca ființă care participă la viața socială. De aceea, setul de indicatori de stare este completat cu o serie de indicatori de evaluare a calității vieții. În anchetele dedicate calității vieții, cum ar fi *Diagnoza Calității Vieții ICCV* sau *European Quality of Life Survey*, indivizii sunt rugați să aprecieze cât de bune sau proaste sunt, de exemplu, serviciile de sănătate publice. De asemenea, sunt rugați să își exprime gradul de mulțumire cu diferite domenii ale vieții proprii, dar și cu viața în general.

Satisfacția cu viața primește o atenție deosebită în studiile de calitate a vieții pentru că reflectă analiza rațională a propriei situații (Diener, 1984), luând în calcul simultan valorile pentru toate criteriile relevante ale standardului subiectiv al unei vieți bune (Veenhoven, 1996). O satisfacție cu viața ridicată înseamnă o calitate a vieții ridicată. Rămâne să identificăm care sunt factorii care sporesc satisfacția cu viața.

Factorii care explică satisfacția cu viața pot fi grupați în mai multe calupuri. Un prim calup se referă la caracteristicile individuale: gen, vârstă, educație, stare civilă, situație financiară, stare de sănătate obiectivă și autoevaluată etc. Acești indicatori sunt nelipsiți, jucând, de regulă, rolul de variabile de control. Un al doilea calup se referă la mecanismele psihologice și psihosociale care determină un nivel mai scăzut sau mai ridicat al satisfacției cu viața, cum ar fi procesul comparației sociale (Michalos, 1985; Easterlin *et al.*, 2010) sau cel al maximizării, întâlnit în

societățile de consum (Schwartz, 2004). Un al treilea calup se poate referi la calitatea percepută a serviciilor publice și condițiilor de trai din zona în care persoana locuiește (Gandelman, Piani și Ferre, 2012). Fără a epuiza subiectul, aș mai nota aici setul de caracteristici materiale și/sau culturale al unei unități sociale cum ar fi vecinătatea (de exemplu, cartierul sau sectorul în oraș), județul, regiunea de dezvoltare sau istorică ori chiar țara (Hagerty și Veenhoven, 2003 ; Inglehart și Welzel, 2005 ; Inglehart *et al.*, 2008 ; Hooghe și Vanhoutte, 2011 ; Mikucka, 2012). Din perspectiva sociologului, abordarea multinivel este necesară pentru explicarea cât mai adecvată a variației satisfacției cu viața. Multinivel înseamnă includerea simultană în analiza de regresie a caracteristicilor individuale culese prin chestionare și a caracteristicilor unității sociale relevante pentru studiu disponibile, de regulă, la Institutele Naționale de Statistică sau Eurostat ori alte organizații internaționale care au ca obiect de activitate agregarea indicatorilor pe care îi furnizează primăriile, spitalele, angajatorii etc. Analiza multinivel este doar o regresie, după cum spune Bickel (2007), dar, fiind ceva mai complicată, nu va fi tratată aici. Pentru a înțelege analiza multinivel, trebuie să înțelegem analiza la nivel individual.

Am putea fi interesați, de exemplu, să vedem cum variază satisfacția cu viața în funcție de starea materială a persoanelor, dar și de reprezentarea despre această stare materială. Teoretic, reprezentarea despre situația materială ar trebui să fie consistentă cu starea concretă. Totuși, așa cum arată multe studii care pornesc de la teoria comparației sociale, starea și reprezentarea pot să acționeze ca factori independenți asupra satisfacției cu viața. De exemplu, unei persoane îi este mai degrabă teamă să piardă un lucru dobândit decât să câștige acel lucru (*loss aversion*) (Tversky și Kahneman, 1991). Tocmai la cei care au acumulat mai multe resurse s-ar putea să acționeze un mecanism de insatisfacție prin modificarea în sus a standardului de referință (Graham și Pettinato, 2006).

Pentru exemplificare vom utiliza datele culese în cercetarea *Diagnoza Calității Vieții 2003* de către Institutul de Cercetare a Calității Vieții. Pentru a menține caracterul introductiv al volumului, vom realiza o analiză de regresie care are doar două variabile independente : venitul persoanelor active pe piața muncii și autopoziționarea pe scala sărac-bogat. Variabila dependentă este satisfacția cu viața. Formulările exacte din chestionar sunt :

Cât de satisfăcut sunteți de viața dvs. în general ?

|                         |   |   |   |   |   |   |   |   |   |                       |
|-------------------------|---|---|---|---|---|---|---|---|---|-----------------------|
| 0                       | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10                    |
| Complet<br>nesatisfăcut |   |   |   |   |   |   |   |   |   | Complet<br>satisfăcut |

|   |               |
|---|---------------|
| Vă rugăm să menționați toate veniturile gospodăriei dvs. din luna trecută, mai 2003 |               |
|   | Dumneavoastră |
| Salariul din activitatea principală   |               |

|   |  |
|---|--|
| Salariu de la un al doilea loc de muncă                       |  |
| Venituri din activități ca întreprinzător/patron              |  |
| Venituri din activități ocazionale                            |  |
| Venituri din proprietăți (profit, dobânzi, dividende, chirii) |  |
| Venituri obținute din vânzarea produselor agricole            |  |

respectiv

În orice societate, unii oameni se consideră bogați, alții se consideră săraci. Având în vedere numerotarea de la 1 la 10, dvs. unde vă situați ?

|       |   |   |   |   |   |   |   |   |       |
|-------|---|---|---|---|---|---|---|---|-------|
| 1     | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10    |
| Sărac |   |   |   |   |   |   |   |   | Bogat |

Regresia liniară multiplă se exprimă formal prin ecuația :

$$Y = \alpha + \beta_1 * X_1 + \beta_2 * X_2 + \dots \beta_n * X_n + e.$$

Pentru noi, aceasta devine :

$$\begin{array}{ccccccc} \text{Satisfacția} & & \text{Venitul} & & \text{Autopозиționarea} & & \\ \text{cu} & = \alpha + \beta_1 * & \text{persoanelor} & + \beta_2 * & \text{pe scala} & & \\ \text{viața} & & \text{active} & & \text{sărac-bogat} & & + e. \end{array}$$

Rulând această analiză putem răspunde următoarelor întrebări :

- Există un efect semnificativ statistic al celor două variabile independente asupra satisfacției cu viața ?
- Dacă există, care este direcția acestei relații ? În ce sens se modifică satisfacția cu viața când venitul persoanelor active se modifică ? Dar când autopозиționarea pe scala sărac-bogat se modifică ?
- Cu câte unități se modifică satisfacția cu viața atunci când venitul persoanelor active se modifică cu o unitate ? Dar când autopозиționarea pe scala sărac-bogat se modifică cu o unitate ?
- Ce parte din varianța satisfacției cu viața este explicată de venitul persoanelor active și autopозиționarea pe scala sărac-bogat ?

În termeni substanțiali, putem afla dacă situația materială și/sau reprezentările despre aceasta explică satisfacția cu viața și, în caz că da, dacă efectul situației materiale se păstrează atunci când controlăm pentru reprezentarea despre aceasta. De asemenea, putem deduce dacă trebuie să mai căutăm și alți factori explicativi ai satisfacției cu viața pe care i-am omis din analiză. Acesta este un exemplu didactic. Dacă am scrie o lucrare științifică, atunci, cu certitudine, modelul ar trebui să fie mai complex. Am exclus, de exemplu, variabile de control esențiale cum ar fi genul, vârsta, educația sau alți predictorii esențiali ai satisfacției cu viața cum ar fi evaluarea propriei stări de sănătate, evaluarea domeniilor importante

ale vieții ș.a. Indiferent de cât de complex este modelul nostru, trebuie, înainte de a începe analiza în SPSS, să avem un model explicativ clar specificat care țină cont de ceea ce a fost deja demonstrat în domeniu. Una dintre cerințele esențiale ale acestei analize este specificarea corectă a modelului, adică includerea tuturor variantelor relevante. Evident, realitatea socială este prea complexă pentru a oferi explicații perfecte. Dar explicațiile parțiale pe care le producem trebuie să fie consistente. De aceea, analiza de regresie nu se face prin „încercare și eroare”. Nu deschidem baza de date și începem să introducem și să scoatem variabile independente în model până când rezultă ceva care seamănă cu ceea ce credeam că ar fi trebuit să rezulte. În fond, în eșantioanele cu volume mari o să găsim relații semnificative statistice din pură întâmplare.

Regresia liniară multiplă este doar un tip de regresie. Probabil, este cel mai utilizat tip. Opțiunea pentru un tip de regresie ține, printre altele, de caracteristicile variabilei dependente. Dacă variabila dependentă este cantitativ continuă, atunci putem utiliza regresia liniară multiplă. Dacă este dummy (1/0), unde codul 1 este atribuit caracteristicii care ne interesează, atunci putem utiliza regresia logistică binară. Dacă este nominală cu cel puțin trei categorii, atunci putem utiliza regresia logistică multinomială. Dacă reprezintă o numărare și are o distribuție în formă de J întors sau J simplu, putem utiliza regresia **count**. Dacă este ordinală, putem utiliza regresia ordinală. Lista poate continua. Literatura în această zonă este bine dezvoltată. Revenind la regresia liniară multiplă, am spus că variabila dependentă trebuie să fie cantitativă continuă. În științele sociale, cu precădere, dar fără a ne limita doar la acestea, este destul de greu să identificăm instrumente de măsurare care produc variabile care iau, teoretic, o infinitate de valori. De regulă, atunci când reușim să măsurăm cantitativ, acestea au un caracter discret. Una dintre cele mai întâlnite proceduri de măsurare în științele sociale este scala tip Likert. Rensis Likert este unul dintre pionierii măsurării în științele sociale, propunând o scală compusă care îi poartă numele fiind, chiar și astăzi, foarte populară (Likert, 1932). Atunci când auzim un analist spunând „scală tip Likert” nu înseamnă în mod necesar că se referă la scala compusă. Acesta s-ar putea referi la tipul variantelor de răspuns. Forma standard este Acord/Dezacord, acestea fiind extremele unei scale de răspuns cu minim patru puncte : Acord total (4), Acord (3), Dezacord (2), Dezacord total (1). Există multe variante, cu sau fără variantă de mijloc :

|   |
|---|
| Dezacord total (1), Dezacord (2), Acord (3), Acord total (4)                                |
| Dezacord total (1), Dezacord (2), Nici acord, nici dezacord (3), Acord (4), Acord total (5) |
| Dezacord total (1) (2) (3) (4) (5) (6) Acord total (7)                                      |
| Dezacord total (1) (2) (3) (4) (5) (6) (7) (8) (9) Acord total (10)                         |
| Dezacord total (1) (2) (3) (4) (5) (6) (7) (8) (9) (10) Acord total (11)                    |
| Complet nesatisfăcut (1) (2) (3) (4) (5) (6) (7) (8) (9) Complet satisfăcut (10)            |

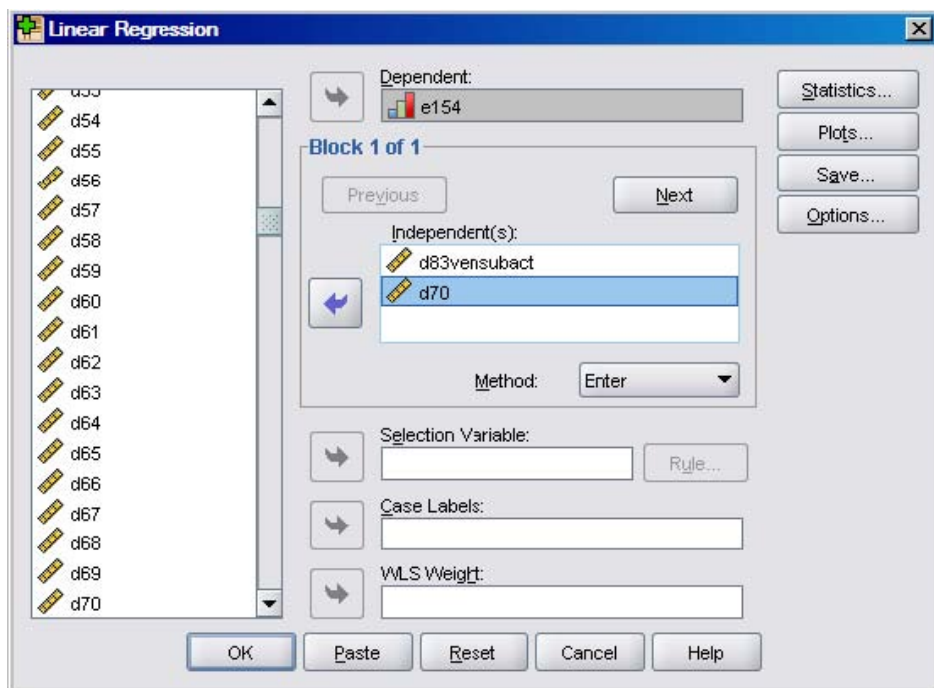
Pe lângă numărul variantelor de răspuns sau opțiunea pentru utilizarea variantei de mijloc, pot fi folosite și etichete diferite. Acestea sunt variabile ordinale care, mai ales când au cel puțin șapte variante de răspuns, sunt considerate de interval. Fiind considerate de interval, sunt utilizate în mod curent în analizele statistice ca variabile dependente în regresia liniară multiplă. Există argumente pro și contra (Carifio și Perla, 2007). Referindu-se mai degrabă la variabilele numerice cu puține valori (cantitative discrete) folosite drept variabile dependente, Berry (1993) recomandă să nu folosim variabilele cantitative discrete ca dependente în regresia liniară atunci când numărul valorilor este mai mic decât 5, iar Fox (1991), pe lângă această recomandare, fără însă a cuantifica ca Berry, consideră că mai reprezintă o problemă serioasă doar atunci când majoritatea răspunsurilor sunt concentrate pe un număr mic de valori.

O altă cerință esențială a regresiei liniare este, așa cum sugerează chiar numele analizei, ca relația dintre variabila dependentă și variabilele independente să fie liniară. Dacă nu este respectată această cerință, atunci trebuie aplicată o formă de regresie nonliniară.

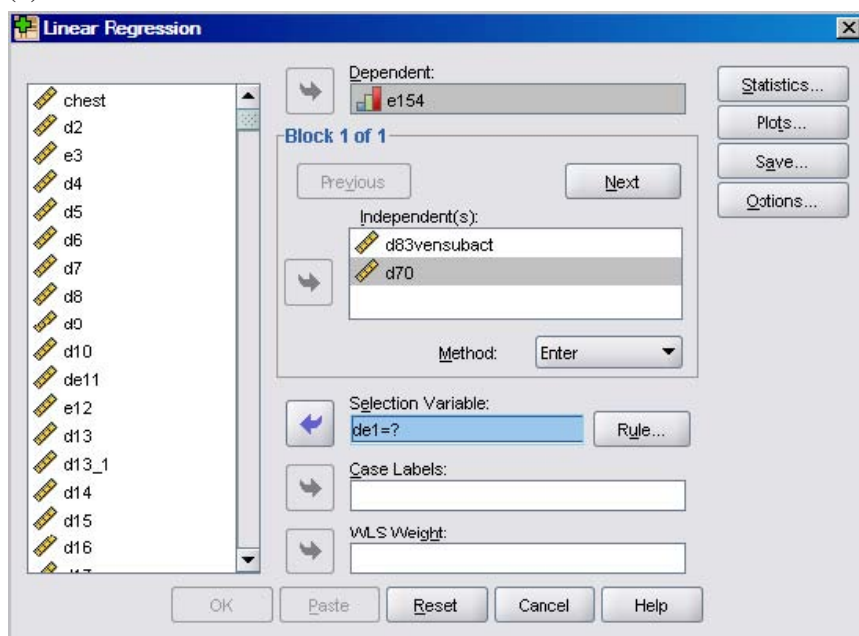
În SPSS, analiza de regresie liniară multiplă poate fi realizată din meniul **Analyze > Regression > Linear** (figura 8.4). Acesta este intuitiv.

**Figura 8.4.** Meniul Analyze > Regression > Linear

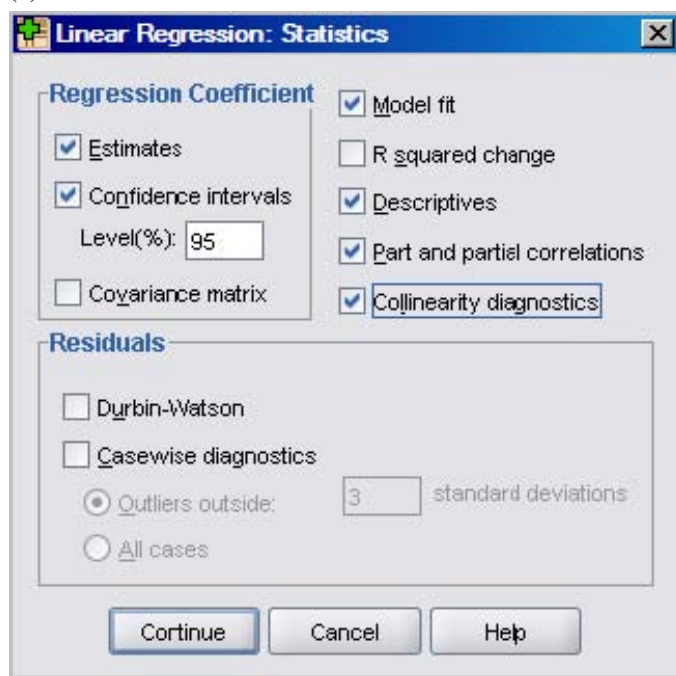
(a)



(b)



(c)



**Linear Regression: Statistics**

**Regression Coefficient**

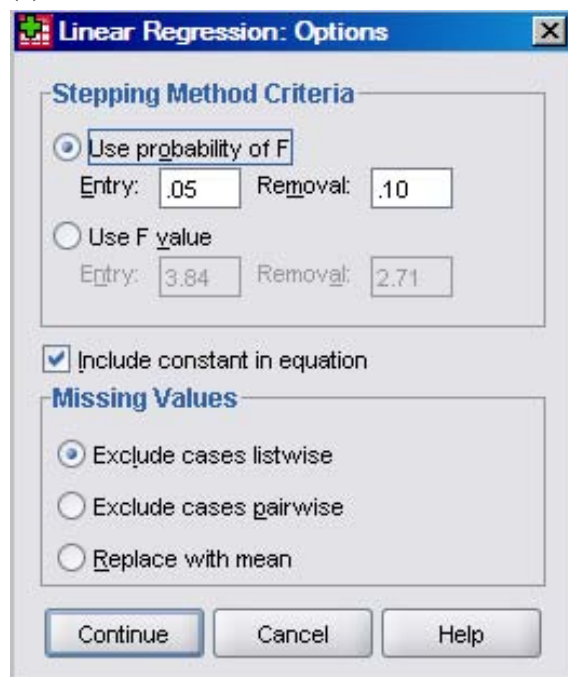
- ☒ Estimates
- ☒ Confidence intervals  
Level(%): 95
- ☐ Covariance matrix
- ☒ Model fit
- ☐ R squared change
- ☒ Descriptives
- ☒ Part and partial correlations
- ☒ Collinearity diagnostics

**Residuals**

- ☐ Durbin-Watson
- ☐ Casewise diagnostics
- ☒ Outliers outside: 3 standard deviations
- ☐ All cases

Continue Cancel Help

(d)



**Linear Regression: Options**

**Stepping Method Criteria**

- ☒ Use probability of F  
Entry: .05 Removal: .10
- ☐ Use F value  
Entry: 3.84 Removal: 2.71

☒ Include constant in equation

**Missing Values**

- ☒ Exclude cases listwise
- ☐ Exclude cases pairwise
- ☐ Replace with mean

Continue Cancel Help



În secțiunea **Dependent**, introducem variabila dependentă. Aici aceasta este satisfacția cu viața care, în baza de date, are numele e154.

Urmează secțiunea **Block 1 of 1 – Independent(s)**, unde introducem variabilele independente. În exemplul nostru, avem două variabile independente, venitul din ultima lună al persoanelor active pe piața muncii, care poartă numele d83vensubact în baza de date, respectiv autopoziționarea pe scala sărac-bogat, care poartă numele d70 în baza de date. Cum spuneam, analiza de regresie presupune elaborarea modelului explicativ a priori utilizării programului de statistică. Putem aborda analiza în mai multe moduri. O primă variantă constă în introducerea tuturor variabilelor independente într-un singur calup sau block (am să folosesc **block** pentru a asigura corespondența cu programul). A doua variantă constă în gruparea, justificată teoretic, a variabilelor independente în mai multe blockuri. Am putea crea un block care conține variabilele de control (gen, vârstă, stare civilă etc.). Apoi am putea crea un alt block care conține informații despre situația materială a persoanei (venit, proprietăți etc.). În fine, am putea crea un block care conține informații despre cum se vede (percepe) persoana în societate din perspectiva resurselor materiale pe care le deține (autopoziționarea pe scala sărac-bogat, raportarea subiectivă a venitului la necesități etc.). Pentru că variabilele din cele trei blockuri au o utilitate proprie, surprinzând aspecte distincte de celelalte, are sens să le folosim ca atare.

Variabilele independente pot fi folosite ca atare sau pot fi grupate în scoruri compuse. Dacă teoria spune că unele variabile ar putea fi grupate sau trebuie grupate în diferite scoruri compozite sau, altfel spus, indici, atunci am fi utilizat în regresie acești indici. De exemplu, dacă am fi măsurat o variabilă independentă printr-o scală compusă Likert, atunci ar fi fost necesară calcularea scorului sumativ (varianta standard) sau am fi calculat media afirmațiilor care o compune ori am fi realizat o analiză factorială exploratorie salvând scorurile factoriale pe care, ulterior, le-am fi utilizat în regresie. Există mai multe metode de calculare a indicilor, decizia aparținând în final analistului. Acesta va trebui să pună în balanță proprietățile statistice ale indicelui calculat cu dificultatea de interpretare a acestuia în analiza de regresie. Scorul sumativ este mai greu de interpretat decât media variabilelor care constituie scala compusă. Indicele calculat ca medie a variabilelor variază în același interval cu cel al variantelor de răspuns, deci va fi mai ușor de înțeles. Scorul sumativ pentru o scală compusă cu 4 variabile și 10 variante de răspuns, unde 1 = acord și 10 = dezacord, variază între 4, dacă respondentul alege codul 1 la toate variabilele, și 40, dacă alege codul 10 la toate variabilele. Cercetătorul trebuie să clarifice ce înseamnă scorul 13 sau scorul 33. La fel se întâmplă cu scorul factorial.

În secțiunea **Method**, avem mai multe metode, cea implicită fiind **Enter**. Aceasta este cea pe care o preferăm deoarece lasă la latitudinea cercetătorului modul în care introduce variabilele independente în analiză. Este consistentă cu elaborarea preliminară a modelului explicativ. Celelalte seamănă, mai degrabă,



cu un proces de încercare-eroare prin care analistul „caută” o relație semnificativă statistic.

În secțiunea **Selection Variable** putem introduce o variabilă care identifică anumite grupuri, instruind astfel SPSS-ul să ruleze analiza de regresie doar pe anumite cazuri. De exemplu, dacă doresc să realizez analiza doar pentru bărbați, atunci introduc variabila sex, aici **de1** (figura 8.4b). Odată introdusă variabila, se activează butonul **Rule**. De la **de1 = ?** trebuie să ajungem la **de1=1**. Adică trebuie să introducem codul care identifică grupul pentru care dorim să facem analiza de regresie. Printr-un tabel de frecvență am aflat că bărbații au codul 1 și, deoarece dorim să rulăm regresia pentru bărbați, apăsăm butonul **Rule** și introducem cifra 1 în câmpul **Value** după ce ne-am asigurat că este selectată opțiunea **equal to** în secțiunea **Define selection rule**.

În secțiunea **Case Labels** putem introduce o variabilă care identifică cazurile în mod precis în graficele pe care le realizăm odată cu celelalte calcule specifice analizei de regresie. De exemplu, am putea introduce identificatorul unic pentru fiecare respondent care, în această bază de date, se numește chest.

În secțiunea **WLS Weight** putem introduce o variabilă specială care ne permite rularea unui regresii liniare ajustate, utilă atunci când este încălcată asumptia homoscedasticității (homoskedasticity) (Lewis-Beck, 1980).

Meniul are o serie de butoane : **Statistics**, **Plots**, **Save** și **Options**. Vom prezenta în continuare unele dintre cele mai importante și utile opțiuni pe care le putem alege.

Butonul **Statistics** (figura 8.4c) conține informațiile esențiale pentru care alegem să rulăm această analiză. Implicit sunt selectate, în secțiunea **Regression Coefficient**, **Estimates** și, alături, **Model fit**. **Estimates** ne va afișa coeficienții de regresie nestandardizați și coeficienții de regresie standardizați. Coeficienții de regresie nestandardizați (acei  $\beta$  din ecuația de regresie) ne arată cu cât se modifică variabila dependentă atunci când variabila independentă corespunzătoare se modifică cu o unitate. Aceștia pot avea semnul minus sau plus, în funcție de relația dintre X și Y, dar și de modul în care sunt codificate cele două. **Model fit** ne afișează valorile  $R^2$  și  $R^2$  ajustat.  $R^2$  se numește coeficient de determinare și ne arată cât din variația variabilei dependente este explicată de variabilele independente incluse în model. Acesta variază între 0 și 1, dar noi îl vom transforma în procente pentru că este mai ușor de citit. Un  $R^2$  egal cu 0.56 înseamnă că 56% din variația variabilei dependente este explicată de variabilele independente incluse în model. Cu cât este mai mare valoarea, cu atât modelul este mai informativ. Lewis-Beck (1980) enumeră următoarele situații pe care trebuie să le avem în vedere când interpretăm valoarea coeficientului de determinare : (a) o valoare mare nu este utilă pentru interpretarea teoretică dacă modelul nu este specificat corect din punct de vedere logic (oferim explicații tautologice) ; (b) o valoare mică nu sugerează în mod necesar un model specificat greșit, această situație putându-se datora unor relații nonliniare între dependentă și independente.

Pentru că  $R^2$  crește odată cu introducerea de noi variabile independente în model, consultăm  $R^2$  ajustat, care ia în calcul această situație. Sunt situații însă, cum ar fi rularea analizei pe eșantioane mici ( $n$  sub 100) folosind multe variabile independente (peste 20), când ajustarea poate da greș (Tabachnick și Fidell, 2007).

În mod uzual, mai bifăm **Confidence intervals, R squared change, Descriptives, Part and partial correlations, Collinearity diagnostics**.

**Confidence intervals** ne oferă intervalele de încredere pentru coeficienții de regresie nestandardizați. Aceștia ne oferă posibilitatea să înțelegem mai realist situația explicativă decât estimarea punctuală. Putem vedea limitele între care poate varia valoarea cu care se modifică variabila dependentă atunci când variabila independentă se modifică cu o unitate. Când intervalul este larg, atunci estimarea nu este tocmai utilă din punct de vedere teoretic (Lewis-Beck, 1980).

**R squared change** este util atunci când utilizăm logica blockurilor. Ne va arăta în ce măsură un nou block de variabile aduce un plus în explicația variabilei dependente. Ca și  $R^2$  ajustat, ia valori între 0 și 1, dar îl citim în procente pentru o interpretare mai ușoară. Cu cât este mai mare valoarea sa, în condițiile unei specificări corecte a modelului explicativ, cu atât contribuția explicativă este mai importantă. Acesta se citește împreună cu valoarea nivelului de semnificație a testului calculat (*sig. F change*): atunci când  $p$  mai mic sau egal cu 0.05, blockul respectiv de variabile contribuie semnificativ statistic la explicarea variabilei dependente.

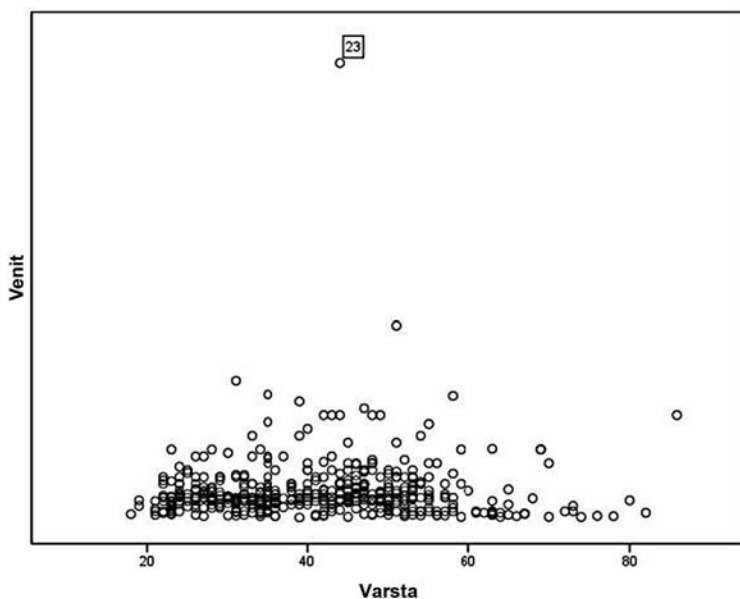
**Descriptives** calculează mediile și abaterile standard pentru fiecare variabilă introdusă în ecuație și ne arată volumul eșantionului pentru care sunt efectuate calculele. Putem considera această opțiune ca un punct de control în analiză. Putem calcula media pentru variabilele introduse în model? Dacă răspunsul este afirmativ, atunci rezultatele analizei de regresie liniare pot fi interpretate. Dacă nu, atunci trebuie să căutăm o soluție pentru variabila unde nu are sens media. Transformarea în variabile dummy (1/0) este soluția atunci când trebuie să utilizăm variabile nominale ca variabile independente. Dacă, de exemplu, trebuie să utilizăm religia ca predictor, aceasta având trei categorii, vom alege o categorie de referință și, cu celelalte două, vom realiza două variabile dummy. Pentru alegerea categoriei de referință nu există o regulă general valabilă: decizia depinde de interesele analistului. De exemplu, dacă religia are categoriile ortodox, catolic și protestant, iar interesul cercetătorului este să compare evoluția dependentei la catolici și protestanți prin raportare la ortodocși, atunci va alege religia ortodoxă ca referință și va crea două dummy-uri astfel:

| Variabila inițială      | Variabila dummy 1:<br>catolic | Variabila dummy 2:<br>protestant |
|-------------------------|-------------------------------|----------------------------------|
| Apartenența religioasă: | 1 devine 0                    | 1 devine 0                       |
| 1. ortodox              | <u>2 devine 1</u>             | 2 devine 0                       |
| 2. catolic              | 3 devine 0                    | <u>3 devine 1</u>                |
| 3. protestant           |                               |                                  |

Este obligatorie introducerea simultană în analiză a celor două variabile dummy. Cele două variabile sunt create folosind meniul **Transform > Recode into different variables**. Media unei variabile dummy indică procentul cazurilor din acea categorie prezente în eșantion și, dacă eșantionul este reprezentativ pentru o populație, respectând structura acesteia, indică procentul cazurilor din acea categorie prezente în populație. De exemplu, dacă media variabilei dummy catolic este egală cu 0.31, atunci avem 31 % catolici în eșantionul nostru.

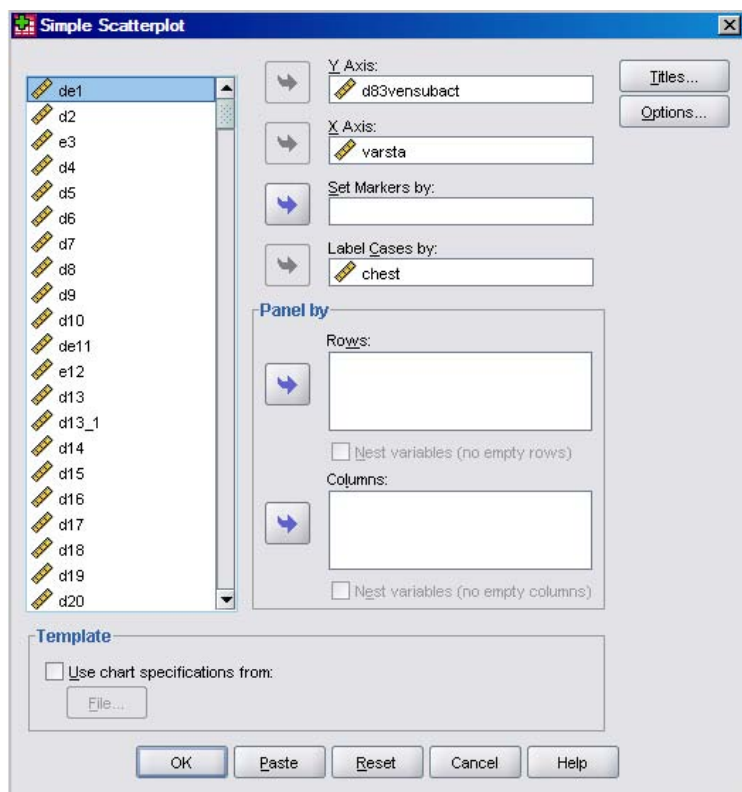
Verificând mediile variabilelor introduse în analiză, ne aducem aminte că acest indicator este influențat de cazurile extreme. Cazurile extreme sunt univariate sau multivariate. Prezența unor astfel de cazuri în analiza de regresie poate modifica serios estimările calculate de program. De aceea, o verificare univariată sau bivariată folosind meniul **Analyze > Descriptives > Explore** este necesară. O altă modalitate complementară, vizuală, constă în realizarea unui grafic scatterplot sau, în limba română, „nor de puncte”. În figura 8.5, care prezintă relația dintre vârsta măsurată în ani împliniți și veniturile persoanelor active pe piața muncii din eșantionul DCV ICCV 2003, observăm un caz extrem : o persoană cu vârsta undeva între 40 și 50 de ani are însumate venituri neașteptat de mari pentru vârsta sa. Probabil acel venit va ieși în evidență și la o inspectare univariată a variabilei respective, dar sunt situații în care nu se întâmplă așa, scatterplotul oferind o informație foarte utilă în acest sens.

**Figura 8.5.** Scatterplot care ne arată un caz extrem



Graficul a fost realizat din meniul **Graphs > Legacy dialogs > Scatter – Dot > Simple Scatterplot** (figura 8.6).

Figura 8.6. Meniul Simple Scatterplot



Considerând vârsta variabila explicativă (logic, nici nu am avea cum să o considerăm altfel) și venitul variabila explicată, cea dintâi este introdusă la **X axis**, iar cea din urmă la **Y axis**. Pentru că, în situația în care observăm vreun caz extrem, dorim să îl identificăm ușor, am introdus la **Label Cases by** variabila care conține id-ul unic al fiecărui respondent care, aici, se numește chest. În figura 8.5 eticheta 23 indică id-ul, această valoare putând fi utilizată pentru filtrarea acestui caz din analizele viitoare, de exemplu.

Aceste informații sunt utile și pentru tabelul care conține corelațiile bivariate, afișat tot prin alegerea opțiunii **Descriptive statistics**. Corelațiile bivariate ne ajută să ne facem o primă idee cu privire la relațiile dintre variabilele incluse în analiză.

Opțiunea **Part and partial correlations** va afișa trei tipuri de corelație: **zero-order**, **part** și **partial**. Dintre acestea ne interesează în mod deosebit corelațiile semiparțiale pe care SPSS le denumește **part correlations**. Această corelație, ridicată la pătrat, ne arată contribuția unică pe care variabila independentă o are la explicarea variabilei dependente. Ne arată cu cât se reduce  $R^2$  dacă acea variabilă independentă este eliminată din ecuația de regresie (Tabachnick și Fidell, 2007). Acești autori explică în detaliu diferența dintre corelația parțială și cea semiparțială,

atrăgând totodată atenția că acest mod de interpretare este specific utilizării regresiei multiple standard, adică cea obținută prin utilizarea metodei **Enter** în SPSS.

În fine, opțiunea **Collinearity Statistics** ne oferă doi indicatori care verifică asumția absenței multicolinearității : **Tolerance** și **VIF (variance inflation factor)**. Schroeder, Sjoquist și Stephan (1986) oferă un exemplu despre ce înseamnă acest lucru : pentru reducerea numărului deceselor rezultate în urma accidentelor auto se introduc simultan două măsuri preventive, purtarea obligatorie a centurii de siguranță și pedepsirea aspră a șoferilor prinși conducând sub influența alcoolului. Deși ambele variabile independente sunt, în esență, importante, va fi greu de distins efectul individual al acestora. Acești autori atrag atenția asupra riscului, atunci când există multicolinearitate, de a întâlni mai des coeficienți nesemnificativi statistic. Exemplul dat de acești autori poate fi completat cu situațiile în care analistul introduce în analiză variabile independente corelate puternic între ele. Corelația puternică poate veni fie din caracterul interșanjabil al indicatorilor (măsoară același lucru), fie din relații de determinare reciprocă. Dacă în exemplul anterior cercetătorul nu poate controla realitatea, interzicând vreuna dintre măsurile preventive, în a doua situație, rolul său este de a analiza anterior analizei de regresie atât din punct de vedere logic, cât și statistic legăturile de determinare dintre variabilele independente. La lista de efecte negative, Field (2009) adaugă și instabilitatea predicției și limitarea valorilor lui  $R^2$ . Revenind la cei doi indicatori, când **Tolerance**, care variază între 0 și 1, are valori mai mici decât 0.1, asumția absenței multicolinearității este încălcată. **VIF** nu are un interval exact de variație. O valoare mai mare decât 10 indică prezența multicolinearității (Field, 2009 ; Kline, 2011).

Revenind la exemplul nostru, să vedem ce se întâmplă cu satisfacția cu viața atunci când controlăm veniturile însumate ale persoanelor active pe piața muncii, respectiv autopozitionarea pe scala sărac-bogat. Ipoteza noastră este că ambele au un efect semnificativ statistic pentru că reflectă mecanisme care nu se suprapun perfect. Resursele materiale ne ajută să ne satisfacem nevoile, dar poziționarea pe scala sărac-bogat implică un proces de comparație socială care ne poate face să ne simțim mai săraci (sau mai bogați) decât suntem. Depinde care este standardul nostru de referință. Înainte de a rula analiza, să apăsăm butonul **Options** (figura 8.4d). Mă opresc asupra secțiunii **Missing Values**, unde este selectat **Exclude cases listwise**. Cunoaștem deja efectele fiecărei metode de tratare a nonrăspunsurilor. Nu o să modificăm nimic în acest meniu.

Outputul analizei este prezentat în continuare. Primul tabel (tabelul 8.4) ne arată media, abaterea standard și volumul eșantionului pentru care este rulată analiza de regresie. Satisfacția cu viața are media egală cu 5.0 și abaterea standard egală cu 2.1. Valorile venitului sunt specifice anului 2003, de aici și modul de prezentare, respectiv sumele mai mici comparativ cu cele de astăzi (atunci când le convertim în lei noi). Observăm că interpretarea mediilor scalelor ordinale, pe care noi le-am considerat de interval, nu este atât de evidentă cum este interpretarea

mediei venitului. Deoarece nonrăspunsurile au fost excluse **listwise**, avem același volum al eșantionului pe care s-a rulat analiza la toate cele trei variabile ( $n = 485$ ). Volumul eșantionului a scăzut semnificativ.

**Tabelul 8.4.** Output regresie liniară multiplă, Descriptives statistics

| Descriptive Statistics                                   |            |                |     |
|--|------------|----------------|-----|
|  | Mean       | Std. Deviation | N   |
| e154 CÂT DE SATISFĂCUT SUNTEȚI DE VIAȚA DVS. ÎN GENERAL? | 5.00       | 2.146          | 485 |
| d83vensubact   | 3973948.45 | 4341032.795    | 485 |
| d70 Poziția pe scala sărăcie-bogăție                     | 4.43       | 1.631          | 485 |

**Tabelul 8.5.** Output regresie liniară, Corelații bivariate

| Correlations        |  |  |              |                                      |
|---------------------|--|--|--------------|--------------------------------------|
|                     |  | e154 CÂT DE SATISFĂCUT SUNTEȚI DE VIAȚA DVS. ÎN GENERAL? | d83vensubact | d70 Poziția pe scala sărăcie-bogăție |
| Pearson Correlation | e154 CÂT DE SATISFĂCUT SUNTEȚI DE VIAȚA DVS. ÎN GENERAL? | 1.000  | .251         | .662                                 |
|                     | d83vensubact   | .251   | 1.000        | .204                                 |
|                     | d70 Poziția pe scala sărăcie-bogăție                     | .662   | .204         | 1.000                                |
| Sig. (1-tailed)     | E154 CÂT DE SATISFĂCUT SUNTEȚI DE VIAȚA DVS. ÎN GENERAL? | .  | .000         | .000                                 |
|                     | D83vensubact   | .000   | .            | .000                                 |
|                     | D70 Poziția pe scala sărăcie-bogăție                     | .000   | .000         | .                                    |
| N                   | E154 CÂT DE SATISFĂCUT SUNTEȚI DE VIAȚA DVS. ÎN GENERAL? | 485  | 485          | 485                                  |
|                     | d83vensubact   | 485  | 485          | 485                                  |
|                     | d70 Pozitia pe scala sărăcie-bogăție                     | 485  | 485          | 485                                  |

Tabelul de corelații din tabelul 8.5 ne oferă o primă imagine a relațiilor care ne interesează, dar și informații preliminare despre asumția absenței multicolinearității. Este calculat coeficientul Pearson, care variază în intervalul  $[-1, 1]$ . Nivelurile de semnificație calculate sunt prezentate în rândul **Sig. (1-tailed)**. În principiu, satisfacția este corelată semnificativ statistic cu ambele variabile. Cei doi indicatori subiectivi, satisfacția și autopoziționarea pe scala sărac-bogat, au o corelație mai puternică, lucru așteptat având în vedere proprietățile psihometrice similare. Corelația mai mică a satisfacției cu venitul nu trebuie luată ca atare pentru că ar putea indica prezența cazurilor extreme sau chiar a unei relații nonlineare.

Tabelul 8.6 prezintă valoarea coeficientului de determinare, **R Square**, și a coeficientului de determinare ajustat, **Adjusted R Square**. Diferențele în acest exemplu sunt mici. De regulă, raportăm ambele valori, oferindu-le posibilitatea cititorilor să aprecieze diferențele. Așadar, 45% din variația satisfacției cu viața pare să fie explicată de venit și autopoziționarea pe scala sărac-bogat. Prima impresie ar fi că am ales bine cele două variabile independente. Testul de semnificație este în tabelul ANOVA. Evident, aceasta este o estimare punctuală, de aceea ne putem imagina că ea poate varia în jurul acestei valori.

**Tabelul 8.6.** Output regresie liniară,  $R^2$

| <b>Model Summary</b>  |                   |          |                   |                            |
|---|-------------------|----------|-------------------|----------------------------|
| Model   | R                 | R Square | Adjusted R Square | Std. Error of the Estimate |
| 1   | .673 <sup>a</sup> | .452     | .450              | 1.591                      |
| a. Predictors: (Constant), d70 Poziția pe scala sărăcie-bogăție, d83vensubact |                   |          |                   |                            |

| <b>ANOVA<sup>b</sup></b>   |            |                |     |             |         |                   |
|--|------------|----------------|-----|-------------|---------|-------------------|
| Model  |            | Sum of Squares | df  | Mean Square | F       | Sig.              |
| 1  | Regression | 1008.351       | 2   | 504.175     | 199.085 | .000 <sup>a</sup> |
|  | Residual   | 1220.647       | 482 | 2.532       |         |                   |
|  | Total      | 2228.998       | 484 |             |         |                   |
| a. Predictors: (Constant), d70 Poziția pe scala sărăcie-bogăție, d83vensubact  |            |                |     |             |         |                   |
| b. Dependent Variable: e154 CÂT DE SATISFĂCUT SUNTEȚI DE VIAȚA DV. ÎN GENERAL? |            |                |     |             |         |                   |

Am ajuns la tabelul (tabelul 8.7) care ne oferă informațiile căutate. Pentru fiecare variabilă independentă ne sunt oferite următoarele informații :

- nivelul de semnificație (coloana **Sig.**) al testului t care indică dacă între variabila independentă și variabila dependentă există o relație semnificativă statistic. Aici, pentru ambele variabile independente, acesta este mai mic decât pragul 0.05, pe care am decis să îl utilizăm ca referință, deci ambele variabile par să influențeze satisfacția cu viața.

- coeficienții de regresie nestandardizați (coloana **Unstandardized Coefficients – B**). Aceștia ne arată că, atunci când venitul crește, satisfacția cu viața crește cu 0.001 puncte pe scală, respectiv că, atunci când individul se consideră mai bogat, satisfacția cu viața crește cu 0.839 puncte pe scală. Valoarea foarte mică a coeficientului venitului poate însemna: (a) efectul venitului este de fapt mic sau inexistent atunci când controlăm pentru autopoziționarea pe scala sărac-bogat; (b) relația dintre venit și satisfacția cu viața nu este liniară, deci ar trebui să revizuim analiza (eliminarea cazurilor extreme dacă există, transformarea variabilelor, introducerea în regresie a pătratului venitului etc.); (c) alte asumptii sunt încălcate.
- coeficienții de regresie standardizați (coloana **Standardized Coefficients – Beta**) sunt folosiți uneori pentru a spune care dintre predictorii are contribuția cea mai importantă la explicarea variabilei dependente. Totuși aceștia nu pot fi interpretați pentru variabilele dummy (Lewis-Beck, 1980), de aceea ne uităm mai degrabă la pătratul corelațiilor semiparțiale din coloana **Correlations – Part**.
- indicii care testează absența multicolinearității sunt prezentați în coloana **Collinearity Statistics**. Indicii de toleranță sunt foarte mari, având valori peste pragul 0.1, iar VIF este mai mic decât 10 pentru ambele independente. Statistic nu există multicolinearitate. Dar trebuie să ne gândim și dacă, logic, efectul individual al celor două variabile independente poate fi disociat.

Tabelul 8.7. Output regresie liniară, Coefficients

| Coefficients <sup>a</sup> |                                      |                             |            |                           |        |      |                                 |             |              |         |      |                         |       |
|---------------------------|--------------------------------------|-----------------------------|------------|---------------------------|--------|------|---------------------------------|-------------|--------------|---------|------|-------------------------|-------|
| Model                     |                                      | Unstandardized Coefficients |            | Standardized Coefficients | t      | Sig. | 95.0% Confidence Interval for B |             | Correlations |         |      | Collinearity Statistics |       |
|                           |                                      | B                           | Std. Error | Beta                      |        |      | Lower Bound                     | Upper Bound | Zero-order   | Partial | Part | Tolerance               | VIF   |
| 1                         | (Constant)                           | 1.046                       | .211       |                           | 4.952  | .000 | .631                            | 1.461       |              |         |      |                         |       |
|                           | d83ven-subact                        | .000                        | .000       | .120                      | 3.499  | .001 | .000                            | .000        | .251         | .157    | .118 | .958                    | 1.043 |
|                           | d70 Poziția pe scala sărăcie-bogăție | .839                        | .045       | .638                      | 18.518 | .000 | .750                            | .928        | .662         | .645    | .624 | .958                    | 1.043 |

a. Dependent Variable: e154 CÂT DE SATISFĂCUT SUNTEȚI DE VIAȚA DVS. ÎN GENERAL?

Am înțeles care este logica regresiei liniare multiple și cum se realizează în SPSS. Pasul următor firesc constă în verificarea tuturor asumptiilor pe care această analiză le are. Pentru înțelegerea lor vă recomand să parcurgeți lucrarea scrisă de Berry (1993).



### 8.3. Exerciții

Pentru aceste exerciții utilizăm baza de date și/sau chestionarul World Values Survey 2012 rezultat(ă/e) în urma aplicării chestionarului în România. Baza de date poate fi descărcată de pe pagina de internet a *Grupului Românesc pentru Studiul Valorilor Sociale* (<http://www.romanianvalues.ro>).

1. Citiți materialul scris de Bogdan Voicu, Horațiu Rusu și Mircea Comșa, cu titlul *Atitudini față de solidaritate în România*, care a fost publicat în volumul coordonat de Lucian Marina, *Ocupare și incluziune socială*, apărut la Editura Presa Universitară Clujeană, în 2013.
2. Creați variabila dependentă „solidaritate”. Înainte de aceasta, rescalați în acord cu modul de lucru al autorilor.
3. Creați variabila „orientare de valoare materialistă sau postmaterialistă” în acord cu modul de lucru al autorilor.
4. Inversați scala care măsoară importanța acordată religiei în acord cu modul de lucru al autorilor.
5. Creați variabila dummy care măsoară comportamentul religios în acord cu modul de lucru al autorilor.
6. Continuați procesul de creare, recodificare, transformare al variabilelor „mândria de a fi român”, „sentimentul apartenenței naționale”, „individualism” și „clasa socială” în acord cu modul de lucru al autorilor.
7. Pregătiți pentru analiză variabilele „vârstă”, „venit”, „educație”, „sex” și „tip de localitate” în acord cu modul de lucru al autorilor.
8. Rulați regresia liniară multiplă în care „solidaritatea” este variabila dependentă, iar toate celelalte sunt independente.
9. Rulați din nou regresia liniară multiplă, dar de data aceasta folosiți blockurile. Ce informație suplimentară obțineți în acest mod?
10. Realizați un raport de două pagini care să descrie rezultatul modelului complet de regresie (cu toți predictorii): pe prima pagină este inserat tabelul de regresie, iar pe a doua pagină acesta este comentat cu trimitere la teoriile din textul celor trei autori.



## Bibliografie

- Agresti, Alan și Finlay, Barbara. 2008. *Statistical methods for the social sciences*. Upper Saddle River, New Jersey : Prentice Hall International, Inc.
- Agresti, Alan și Franklin, Christine. 2013. *Statistics. The art and science of learning from data*. Boston : Pearson.
- Berry, William D. 1993. *Understanding Regression Assumptions*. Newbury Park, CA : Sage.
- Bickel, Robert. 2007. *Multilevel Analysis for Applied Research. It's just a regression !* New York : The Guilford Place.
- Bradburn, Norman M., Sudman, Seymour și Wansink, Brian. 2004. *Asking Questions. The Definitive Guide to Questionnaire Design — For Market Research, Political Polls, and Social and Health Questionnaires*, Revised Edition. San Francisco : Jossey-Bass, A Wiley Imprint.
- Carifio, James și Perla Rocco, J. 2007. Ten common misunderstandings, misconceptions, persistent myths and urban legends about Likert scales and Likert response format and their antidotes. *Journal of Social Sciences* 3(3) : 106-16.
- Carmines, Edward G. și Zeller, Richard A. 1979. *Reliability and validity assesement*. London : Sage.
- Carroll, John B. 1961. The nature of the data, or how to choose a correlation coefficient. *Psychometrika* 26(4).
- Chambers, John M., Cleveland, William S., Kleiner, Beat și Tukey, Paul A. 1983. *Graphical methods for data anaylsis*. Pacific Groove, California : Wadsworth & Brooks / Cole Publishing Company.
- Chelcea, Septimiu. 2007. *Metodologia cercetării sociologice. Metode cantitative și calitative*. București, Editura Economică.
- Chen, Peter Y. și Popovich, Paula M. 2002. *Correlation. Parametric and Nonparametric Measures*. Thousand Oaks, CA : SAGE.
- Cramer, Duncan și Howitt, Dennis. 2004. *The SAGE Dictionary of Statistics. A Practical Resource for Students in The Social Sciences*. London, Thousands Oaks : SAGE.
- Cummins, Robert A. 2003. Normative life satisfaction: Measurement issues and a homeostatic model. *Social Indicators Research* 64(2) : 225-56.
- de Vaus, David. 2002. *Analyzing social science data. 50 key problems in data analysis*. London : SAGE.
- Diener, Ed. 1984. Subjective Well-Being. *Psychological Bulletin* 95(3) : 542-75.
- Easterlin, Richard A., Angelescu McVey, Laura, Switek, Małgorzata, Sawangfa Onnichia și Smith Zweig, Jakueline. 2010. The happiness-income paradox revisited. *Proceedings for the National Academy of Science of the United States of America* 107(52) : 22463-68.
- Field, Andy. 2009. *Discovering statistics using SPSS (ans sex and drugs and rock 'n' roll)*. Los Angeles : SAGE.

- Fox, John. 1991. *Regression diagnostics*. Newbury Park, CA : SAGE.
- Fox, John. 2009. *A mathematical primer for social statistics*. Los Angeles : SAGE.
- Gandelman, Nestor, Piani, Giorgia și Ferre, Zuleika. 2012. Neighborhood Determinants of Quality of Life. *Journal of Happiness Studies* 13 : 547-63.
- Good, Phillip I. și Hardin, James W. 2012. *Common errors in statistics (and how to avoid them)*. Hoboken, New Jersey : John Wiley and Sons, Inc.
- Graham, Carol și Pettinato, Stefano. 2006. Frustrated achievers : winners, losers, and subjective well-being in Peru's emerging economy. *The ANNALS of the American Academy of Political and Social Science* 606 : 128-53.
- Hagerty, Michael R. și Veenhoven, Ruut. 2003. Wealth and happiness revisited. Growing wealth of nations does go with greater happiness. *Social Indicators Research* 64 : 1-27.
- Hair, Joseph F., Black, William C., Babin, Barry J. și Anderson, Rolph E. 2010. *Multivariate data analysis. A global perspective*. Upper Saddle River N.J. : Pearson.
- Hartwig, Frederick și Dearing, Brian E. 1979. *Exploratory data analysis*. Newbury Park, CA : SAGE.
- Henkel, Ramon E. 1976. *Tests of significance*. Newbury Park, CA : SAGE.
- Hooghe, Marc și Vanhoutte, Bram. 2011. Subjective Well-Being and Social Capital in Belgian Communities. The Impact of Community Characteristics on Subjective Well-Being Indicators in Belgium. *Social Indicators Research* 100(1) : 17-36.
- Inglehart, Ronald, Foa, Roberto, Peterson, Christopher și Welzel, Christian. 2008. Development, freedom, rising happiness. A global perspective, 1981-2007. *Perspectives on Psychological Science* 264-285(3).
- Inglehart, Ronald și Welzel, Christian. 2005. *Cultural Change and Democracy : The Human Development Sequence*. New York and Cambridge : Cambridge University Press.
- Iversen, Gudmund R. și Norpoth, Helmut. 1987. *Analysis of variance*. Thousand Oaks, CA : SAGE.
- Jaccard, James și Jacoby, Jacoby. 2010. *Theory Construction and Model-Building Skills. A Practical Guide for Social Scientists*. New York : The Guilford Press.
- Jacoby, William G. 1997. *Statistical graphics for univariate and bivariate data*. Newbury Park, CA : Sage.
- Jacoby, William G. 1998. *Statistical graphics for visualizing multivariate data*. Newbury Park, CA : Sage.
- Kline, Rex B. 2004. *Beyond significance testing. Reforming data analysis methods in behavioral research*. Washington, DC : American Psychological Association.
- Kline, Rex B. 2011. *Principles and Practice of Structural Equation Modeling*. New York : The Guilford Press.
- Lelkes, Orsolya. 2008. *Happiness Across the Life Cycle : Exploring Age-Specific Preferences*. Policy Brief (2).
- Levy, Paul S. și Lemeshow, Stanley. 2008. *Sampling of populations. Methods and applications*. Hoboken, New Jersey : John Wiley & Sons, Inc.
- Lewis-Beck, Michael S. 1980. *Applied Regression : An Introduction*. Newbury Park, CA : SAGE.
- Likert, Rensis. 1932. *A technique for the measurement of attitudes*. *Archives of Psychology* 22(140) : 5-55.
- Malhotra, Naresh K. și Birks, David F. 2007. *Marketing Research. An Applied Approach*. Harlow, England : Prentice Hall Financial Times.

- Mărginean, Ioan. 1982. *Măsurarea în sociologie*. București : Editura Științifică și Enciclopedică.
- Mărginean, Ioan. 2005. *Semnificația cercetărilor de calitate a vieții*. Pp. 25-60 în *Calitatea vieții în România*, editată de Ioan Mărginean, Ana Bălașa și Cătălin Zamfir. București : Editura Expert.
- Michalos, Alex C. 2005. *Multiple discrepancies theory (MDT) (1985)*. Pp. 305-72 in *Citation Classics from Social Indicators Research. The Most Cited Articles Edited and Introduced by Alex C. Michalos, edited by Alex C. Michalos*. Dordrecht, The Netherlands : Springer.
- Mikucka, Małgorzata. 2012. *Individualist culture lowers well-being of the unemployed due to weaker family support norm. Evidence for Europe*. CEPS/Instead, Luxembourg.
- Mohr, Lawrence B. 1990. *Understanding Significance Testing*. Newbury Park, CA : SAGE.
- Raudenbush, Stephen W., Bryk, Anthony S., Cheong, Yuk Fail, Congdon, Richard T. Jr. și du Thoit, Mathilda. 2011. *HLM 7. Hierarchical Linear and Nonlinear Modeling*. Lincolnwood, IL. : SSI Scientific Software International.
- Reynolds, H.T. 1984. *Analysis of nominal data*. Newbury Park, CA : Sage.
- Rotariu, Traian (coord.), Bădescu, Gabriel, Culic, Irina, Mezei, Elemer și Mureșan, Cornelia. 2006. *Metode statistice aplicate în științele sociale*. Iași : Polirom.
- Rughiniș, Cosima. 2007. *Explicația sociologică*. Iași : Polirom.
- Saris, Willem E. și Gallhofer, Irmtraud N. 2007. *Design, evaluation, and analysis of questionnaires for survey research*. Hoboken, New Jersey : John Wiley & Sons, Inc.
- Schroeder, Larry D., Sioquist, David L. și Stephen, Paula E. 1986. *Understanding regression analysis. An introductory guide*. Newbury Park, CA : SAGE.
- Schwartz, Barry. 2004. *The paradox of choice. Why more is less*. New York : Harper Collins.
- Stoop, Ineke A.L. 2005. *The Hunt for the Last Respondent. Nonresponse in sample surveys*. Social and Cultural Planning Office for Netherlands : The Hague.
- Tabachnick, Barbara G. și Fidell, Linda S. 2007. *Using Multivariate Statistics*. Boston : Pearson.
- Tufte, Edward R. 2001. *The visual display of quantitative information*. Cheshire, Connecticut : Graphics Press.
- Tversky, Amos și Kahneman, Daniel. 1991. *Loss Aversion in Riskless Choice : A Reference-Dependent Model*. *The Quarterly Journal of Economics* 106 (4 [Nov.]) : 1039-61.
- Veenhoven, Ruut. 1996. The study of life satisfaction. Pp. 11-48 in *A comparative study of satisfaction with life in Europe*, editată de Saris, W.E., Veenhoven, R., Scherpenzeel, A.C. și Bunting, B : Eötvös University Press.
- Voicu, Bogdan, Rusu, Horațiu și Comșa, Mircea. 2013. Atitudini față de solidaritate în România. Pp. 17-44 în *Ocupare și incluziune socială*, Lucian Marina. Cluj : Presa Universitară Clujeană.
- Zamfir, Cătălin (coordonator), Popescu, Ion-Andrei, Ștefănescu, Ștefan, Teodorescu, Alin, Vlăsceanu, Lazăr și Zamfir, Elena. 1984. *Indicatori și surse de variație a calității vieții*. București : Editura Academiei Republicii Socialiste România.
- Zimmerman, Anke C. și Easterlin, Richard A. 2006. Happily Ever After? Cohabitation, Marriage, Divorce, and Happiness in Germany. *Population and Development Review* 32(3) : 511-28.



# COLLEGIUM

## Sociologie. Antropologie

au mai apărut:

- Petru Iluț – *Sinele și cunoașterea lui. Teme actuale de psihosociologie*  
Marie-Odile Géraud, Olivier Leservoisier, Richard Pottier – *Noțiunile-cheie ale etnologiei. Analize și texte*  
Marian Preda – *Politica socială românească între sărăcie și globalizare*  
Emile Durkheim – *Regulile metodei sociologice*  
Albert Ogien – *Sociologia devianței*  
Traian Rotariu – *Demografie și sociologia populației. Fenomene demografice*  
Ioan Mihăilescu – *Sociologie generală. Concepte fundamentale și studii de caz*  
W. Richard Scott – *Instituții și organizații*  
Irina Culic – *Metode avansate în cercetarea socială. Analiza multivariată de interdependență*  
Petru Iluț – *Valori, atitudini și comportamente sociale. Teme actuale de psihosociologie*  
Cătălin Zamfir – *O analiză critică a tranziției. Ce va fi „după”*  
Gilles Ferréol, Guy Jucquois (coord.) – *Dicționarul alterității și al relațiilor interculturale*  
Robert K. Yin – *Studiul de caz. Designul, analiza și colectarea datelor*  
Richard A. Krueger, Mary Anne Casey – *Metoda focus grup. Ghid practic pentru cercetarea aplicată*  
Ronald F. King – *Strategia cercetării. Treisprezece cursuri despre elementele științelor sociale*  
Petru Iluț – *Sociopsihologia și antropologia familiei*  
Dumitru Sandu – *Dezvoltare comunitară. Cercetare, practică, ideologie*  
Cătălin Zamfir – *Spre o paradigmă a gândirii sociologice*  
Mircea Agabrian – *Analiza de conținut*  
Adrian Hatos – *Sociologia educației*  
Cătălin Zamfir, Laura Stoica (coord.) – *O nouă provocare : dezvoltarea socială*  
Marian Preda – *Comportament organizațional*  
Mihai Păunescu – *Organizare și câmpuri organizaționale. O analiză instituțională*  
Robert Atkinson – *Povestea vieții. Interviu*  
Traian Rotariu, Petru Iluț – *Ancheta sociologică și sondajul de opinie. Teorie și practică (ediția a II-a)*  
Amia Lieblich, Rivka Tuval-Mashiach, Tamar Zilber – *Cercetarea narativă. Citire, analiză și interpretare*  
Lazăr Vlăsceanu – *Sociologie și modernitate. Tranziții spre modernitatea reflexivă*  
Jacques Coenen-Huther – *Sociologia elitelor*  
Cosima Rughiniș – *Explicația sociologică*  
Dumitru Sandu (coord.), Cosmin Câmpean, Lucian Marina, Mihaela Peter, Vasile Șoflău – *Practica dezvoltării comunitare*  
Cătălin Zamfir, Simona Stănescu (coord.) – *Enciclopedia dezvoltării sociale*  
Mihai Coman – *Introducere în antropologia culturală. Mitul și ritul*  
Traian Rotariu – *Demografie și sociologia populației. Structuri și procese demografice*  
Vintilă Mihăilescu – *Antropologie. Cinci introduceri (ediția a II-a)*  
Remus Gabriel Anghel, István Horváth (coord.) – *Sociologia migrației. Teorii și studii de caz românești*  
Nicu Gavriluță – *Antropologie socială și culturală*  
Petru Iluț – *Psihologie socială și sociopsihologie. Teme recurente și noi viziuni*  
Raluca Popescu – *Introducere în sociologia familiei. Familia românească în societatea contemporană*  
Călin Cotoi – *Introducere în antropologia politică*  
Marian Preda (coord.) – *Riscuri și inechități sociale în România. Raportul Comisiei Prezidențiale pentru Analiza Riscurilor Sociale și Demografice*  
Vintilă Mihăilescu (coord.) – *Etnografia urbane. Cotidianul văzut de aproape*  
Earl Babbie – *Practica cercetării sociale*  
Victor Jupp (coord.) – *Dicționar al metodelor de cercetare socială*

Cristina Gavriluță, Nicu Gavriluță – *Sociologia sportului. Teorie, metode, aplicații*  
 Dumitru Sandu – *Lumile sociale ale migrației românești în străinătate*  
 Traian Rotariu – *Studii demografice*  
 Mihaela Vlăsceanu – *Economie socială și antreprenoriat. O analiză a sectorului nonprofit*  
 Laura Grünberg (coord.) – *Introducere în sociologia corpului. Teme, perspective și experiențe întrupate*  
 Lazăr Vlăsceanu (coord.) – *Sociologie*  
 Alex Preda – *Introducere în sociologia piețelor. Informație, cunoaștere și viață economică*  
 Martine Segalen – *Sociologia familiei*  
 Camelia Beciu – *Sociologia comunicării și a spațiului public. Concepte, teme, analize*  
 Mihai Dinu Gheorghiu, Monique de Saint Martin (coord.) în colaborare cu Bénédicte de Montvalon – *Educație și frontiere sociale : Franța, România, Brazilia, Suedia*  
 Gabriel Jderu – *Introducere în sociologia emoțiilor*  
 Traian Rotariu, Vergil Voineagu (coord.) – *Inerție și schimbare. Dimensiuni sociale ale tranziției în România*  
 Cosima Rughiniș – *Măsurarea sociologică. Teorii și practici ale cuantificării*  
 Lazăr Vlăsceanu – *Introducere în metodologia cercetării sociologice*  
 Petru Iluț (coord.) – *În căutare de principii. Epistemologie și metodologie socială aplicată*  
 Nicu Gavriluță – *Sociologia religiilor. Credințe, ritualuri, ideologii*  
 Marian-Gabriel Hâncean – *Rețelele sociale. Teorie, metodologie și aplicații*  
 Marian Vasile – *Introducere în SPSS pentru cercetarea socială și de piață. O perspectivă aplicată*

[www.polirom.ro](http://www.polirom.ro)

Redactor: Dan Mironescu  
 Coperta: Radu Răileanu  
 Tehnoredactor: Radu Căpraru

Bun de tipar: februarie 2014. Apărut: 2014  
 Editura Polirom, B-dul Carol I nr. 4 • P.O. BOX 266  
 700506, Iași, Tel. & Fax: (0232) 21.41.00; (0232) 21.41.11;  
 (0232) 21.74.40 (difuzare); E-mail: [office@polirom.ro](mailto:office@polirom.ro)  
 București, Splaiul Unirii nr. 6, bl. B3A, sc. 1, et. 1,  
 sector 4, 040031, O.P. 53 • C.P. 15-728  
 Tel.: (021) 313.89.78; E-mail: [office.bucuresti@polirom.ro](mailto:office.bucuresti@polirom.ro)

---

Tiparul executat la S.C. Tipo-Lidana S.R.L.  
 Calea Unirii nr. 35, Suceava  
 Tel.: 0230/517.518; Fax: 0330/401.062  
 E-mail: [office@tipolidana.ro](mailto:office@tipolidana.ro); [www.tipolidana.ro](http://www.tipolidana.ro)

---